

مروری بر روش‌های مبتنی بر درخت تصمیم^۱ جهت پیش‌بینی بیماری دیابت

سعید محمدی *، محمد حسین ندیمی

گروه مهندسی کامپیوتر، دانشکده کامپیوتر، دانشگاه آزاد اسلامی، نجف آباد، اصفهان - ایران

خلاصه

افزایش تعداد بیماران دیابتی، عواقب این بیماری، هزینه‌ها و آمار مرگ و میر مرتبط با آن تلاش برای توسعه و بهبود سیستمهای کامپیوتری تشخیص و پیش‌بینی بر پایه داده کاوی را معنا دار می‌کند. داده کاوی پژوهشی می‌تواند کمک بزرگی به جامعه پژوهشی در کشف الگوهای پنهان، روابط معنا دار و مهم از بین حجم انبوهی از داده‌ها و اطلاعات موجود با توجه رشد جمعیت بیماران باشد و در نهایت به افزایش صحت پیش‌بینی بیماری، سادگی روند درمان و کاهش هزینه‌های مربوط به آن کمک کند. هدف این مطالعه مقایسه تعدادی از کارهای انجام شده اخیر در رابطه با پیش‌بینی و تشخیص بیماری دیابت بوسیله درخت تصمیم است. در این مقاله روش‌های ساخت درختان تصمیم بر پایه نوع شکاف داده و انتخاب ویژگی و کشف رابطه آن با صحت نهایی مدل پیش‌بینی با توجه به تعداد رکوردهای مجموعه داده، متغیرهای ورودی خروجی و ویژگی‌های انتخاب شده توسط شخص پژوهش مورد بررسی قرار می‌گیرد. در این مطالعه روش‌های مختلف شکاف داده با مجموعه داده‌های مختلف از بیماران دیابتی و غیر دیابتی و فاکتورهای ریسک^۲ متفاوت به همراه صحت نهایی هر کدام از مدل‌های پیش‌بینی بررسی و ارزیابی گردید. نتایج تحلیلی نشان داد که تعداد نمونه‌های مجموعه داده و روش‌های شکاف داده در میزان صحت طبقه بند، با هم در ارتباطند و همچنین انتخاب تعداد زیادی از فاکتورهای ریسک بخاطر عدم ارتباط برخی از آنها با تشخیص بیماری در افزایش صحت پیش‌بینی بی تاثیر است.

واژه‌های کلیدی: داده کاوی پژوهشی، پیش‌بینی، درخت تصمیم، بیماری دیابت

۱. مقدمه

امروزه بهمراه رشد روز افزون جمعیت، کاهش فعالیت بدنی افراد به علت گسترش زندگی ماشینی، مصرف غذایی که برای سلامتی مضر هستند و همچنین روش‌های متفاوت زندگی، علاوه بر این که سن ابتلا به دیابت کاهش یافته است، نرخ ابتلا به آن افزایش یافته است. از طرفی بیماری دیابت می‌تواند به خطراتی همچون نایینایی و سکته قلبی منجر شود. هزینه جهانی مصرف شده در رابطه با این بیماری در سال ۲۰۱۳ برابر ۵۴۸ میلیون دلار می‌باشد. تعداد بیماران دیابتی در سال ۲۰۱۴

¹ Decision Tree

* Corresponding author: Email: saeed-mohamadi@sco.iaun.ac.ir

² Risk Factors

برابر ۳۸۷ میلیون نفر بوده است و پیش‌بینی شده است تا سال ۲۰۴۰ حدود ۶۴۲ میلیون نفر در سراسر جهان مبتلا به بیماری دیابت باشند [۱].

در دهه‌های اخیر حجم اطلاعات و داده‌های بیماران بسیار گسترده شده است و شخص پزشک نمی‌تواند روی تمامی فاکتورهای موثر آن تحلیل و ملاحظات مربوطه را داشته باشد؛ همچنین فاکتورهای ریسک متفاوتی در ارتباط با بیماری دیابت وجود دارد که هر کدام تاثیرات متفاوتی در فرایند تشخیص و پیش‌بینی این بیماری دارند. بنابراین استفاده از سیستم‌های تحلیلی کامپیوتری با توجه به آمار افزایشی ابتلا به این بیماری می‌تواند تاثیر زیادی در کشف الگوهای مربوط به تشخیص و پیش‌بینی دیابت و در نهایت کاهش هزینه‌های درمان آن داشته باشد. داده کاوی یک جستجوی خودکار روی منابع بسیار بزرگ از داده‌ها برای یافتن الگوها و رفتارهایی که منجر به کشف دانش می‌شوند است [۲]. یکی از تکنیک‌های مشهور ساده و قدرتمند داده کاوی که در زمینه‌های مختلف از جمله پزشکی استفاده می‌شود، طبقه بند درخت تصمیم است. شهرت این تکنیک به خاطر پیاده سازی ساده آن برای حجم بالایی از داده در کنار صحت بالا و راحتی درک و تفسیر آن توسعه عموم مردم است [۳]. همین موضوع دلیل ما برای انتخاب طبقه بند درخت تصمیم در کار مقایسه‌ای فعلی جهت پیش‌بینی بیماری دیابت می‌باشد.

تا کنون مطالعات زیادی در مورد پیش‌بینی بیماری دیابت بوسیله درختان تصمیم انجام شده است. برای مثال در سال ۲۰۱۳ آقای Xue-Hui.m و همکارانش بوسیله درخت تصمیم C5.0 و یک مجموعه داده با ۱۸۴۷ نمونه از افراد دیابتی و غیر دیابتی یک مدل پیش‌بینی را تولید کرده‌اند و به صحت ۷۷.۸۷٪^۱ که بیشترین مقدار در مقایسه با دو طبقه بند دیگر در همان مطالعه بود دستیافتنند [۴]. همچنین در سال ۲۰۱۴ آقای O.Purnik و همکارانش بوسیله درخت^۲ CART و مجموعه داده با ۶۶۴۷ نفر افراد غیر دیابتی توانستند با صحت ۹۰.۵۰٪ یک مدل تشخیص و پیش‌بینی دیابت را بسازند که در ۱۱٪ از افراد این بیماری پیش‌بینی شد. در مطالعات بیان شده و کارهای انجام شده اخیر به مقایسه درخت‌های تصمیم مختلف با توجه به تعداد رکوردهای مجموعه داده اشاره نشده است. اکثر مطالعات از چندین تکنیک طبقه بندی مختلف یا طبقه بند مجموعه ترکیبی^۳ جهت بکارگیری چند تکنیک طبقه بندی در جهت ساخت مدل پیش‌بینی استفاده کرده‌اند.

درختان تصمیم را می‌توان با توجه به روش‌های شکاف داده و انتخاب و یزگی به دسته‌های زیر تقسیم کرد:

- شکاف بر پایه مفهوم Entropy و Information Gain
- شکاف بر پایه Ratio تقسیم می‌شود [۵].
- شکاف بر پایه روش Gini Index
- شکاف بر پایه روش آماری Chi-Squared Test

¹ Classification And Regression Tree

² Ensemble Classifier

شکاف بر پایه آزمایش Chi-Squared^۵ در درخت تصمیم CHAID استفاده شده است که اغلب برای کاربردهای آماری و بازاریابی استفاده می‌شود و در حوزه پزشکی و کارهای مقایسه شده در این مطالعه کاربردی نداشته است. هدف این مطالعه تحلیل ۷ مورد از کارهای انجام شده اخیر در زمینه پیش‌بینی بیماری دیابت بوسیله درختان تصمیم با رویکرد مقایسه آنها از نظر الگوریتم‌های شکاف داده و انتخاب آن با تعداد رکوردهای مجموعه داده، ویژگی‌ها و فاکتورهای ریسک می‌باشد. مقایسه تحلیلی کارهای انجام شده با توجه به مقایسه تکنیک‌های شکاف داده می‌تواند کمکی به کشف رابطه آن با تعداد رکوردهای مجموعه داده و فاکتورهای موثر برای تشخیص و پیش‌بینی این بیماری شود و درنهایت افزایش صحت طبقه‌بندی و کاهش هزینه‌های زمانی ساخت مدل‌های پیش‌بینی را در ارتباط با حذف فاکتورهای ریسک کم ارزش یا نامرتبط داشته باشیم. نتایج مقایسه ای نشان داد که با تعداد تعداد رکوردهای کم و یکسان، روش Information Gain صحت Information Gain Ratio بسته به Gain ارتباط مستقیمی با صحت پیش‌بینی دارد. با مقایسه روش Gini Index و Information Gain به این نتیجه رسیدیم که افزایش تعداد ویژگی‌های فاکتورهای ریسک بیماری در روند ساخت مدل پیش‌بینی به خاطر تفاوت میزان تاثیر آنها با این بیماری نمی‌تواند ارتباط مستقیمی با افزایش صحت مدل داشته باشد. در ادامه این مقاله به تعاریف و آمارهای مربوط به بیماری دیابت در بخش ۲ اشاره می‌شود؛ در قسمت ۳ به مفاهیم مربوط به داده کاوی پزشکی و انواع درختان تصمیم می‌پردازیم بعد از آن در قسمت ۴ توضیح مختصری درمورد تکنیک‌های طبقه‌بندی استفاده شده در حوزه داده کاوی و سپس روش ارزیابی طبقه‌بندی‌ها داده می‌شود و در قسمتهای ۵ و ۶ به ترتیب، مقایسه تحلیلی کارهای اخیر انجام شده در زمینه پیش‌بینی بیماری دیابت بوسیله درخت تصمیم و پیش‌بینی دیابت^۶ بیان می‌شود.

۲. بیماری دیابت

دیابت یک بیماری مزمن یا یک وضعیت متابولیکی (سوخت و ساز) در بدن است که با بالارفتن سطح قند خون در بدن رخ می‌دهد. بالا رفتن قند خون از اینجا ناشی می‌شود که لوزالمudedه هیچ انسولینی تولید نکند، انسولین تولید شده ناکافی باشد و یا انسولین بخوبی کار خود را انجام ندهد [۶]. هورمون انسولین توسط لوزالمudedه تولید و باعث شده تا قند خون وارد سلولهای بدن شود و سوخت و ساز صورت گیرد. در بدن انسان از قند خون به عنوان سوخت جهت تولید انرژی استفاده می‌شود [۷]. دیابت می‌تواند منجر به تخرب سلول‌های عصبی و رگ‌های خون، نابینایی، نارسایی و از کار افتادگی کلیه، قطع عضو و افزایش ریسک سکته و بیماریهای قلبی شود [۸، ۹]. بر طبق گفته ADA^۶ سطح قند خون ناشتا بین ۱۰۰ و ۱۲۵ میلی گرم در دسی لیتر به عنوان پیش دیابت در نظر گرفته می‌شود [۷].

۱.۲. انواع دیابت

بیماری دیابت انواع مختلفی دارد که به ترتیب در مورد آنها توضیحاتی خواهیم داد.

^۵ Chi-squared Automatic Interaction Detector

^۶ American Diabetes Association

۱. دیابت نوع ۱

اگر سلول های تولید انسولین در لوزالمعده(که سلول های بتا^۷ نام دارند) از بین رفته باشند ، هیچ انسولینی تولید نشده و میزان قند خون افزایش می یابد. این نوع دیابت در هرسنی می تواند رخ دهد اما معمولاً در دوران کودکی و یا تا قبل از ۴۰ سالگی پدیدار می شود. علائم و نشانه های آن معمولاً بسیار سریع به طور معمول بعد از چند هفته آشکار می شوند[۷، ۱۰]. دیابت نوع یک ۵ تا ۱۰٪ از تعداد مبتلایان به این بیماری را شامل می شود[۹]. چندین برنامه پیش بینی بصورت آزمایشی در کلینیک های مختلف جهان با مطالعات اضافی برنامه ریزی شده در حال انجام است ولی تا کنون هیچ راه شناخته شده ای برای جلو گیری از دیابت نوع ۱ وجود نداشته است[۱۰].

۲. دیابت نوع ۲

اگر لوزالمعده بتواند مقداری انسولین تولید کند اما نه به اندازه کافی و یا انسولین تولید شده به درستی کار نکند ، این نوع دیابت رخ می دهد. اغلب از سن ۲۵ سالگی شروع شده و در افراد بالای ۴۰ سال نیز مشاهده می شود. همچنین دیابت نوع ۲ بین کودکان، نوجوانان و جوانان به صورت فزاینده ای در حال افزایش است. علائم و نشانه های ممکن است به راحتی واضح و آشکار نباشد زیرا شرایط آن به آرامی در طی یک دوره یکساله گسترش می یابد و ممکن است تنها در یک آزمایش پزشکی روتین مشخص شود[۷]. در این نوع از دیابت ، سلول های بتا در لوزالمعده توانایی خود در تولید انسولین کافی و موثر را رفته رفته از دست می دهند[۱۰]. دیابتی های نوع ۲ نیاز به داروی انسولین برای زنده ماندن ندارند ولی ۲۰٪ از آنها برای کنترل سطوح قند خون با انسولین درمان می شوند. دیابت نوع دو حدود ۹۰٪ از تعداد مبتلایان به این بیماری را شامل می شود و تشخیص آن در مراحل اولیه کاملاً یک کار چالش برانگیز به دنبال وابستگی پیچیده فاکتور های مختلف است [۱۱].

۳. پیش دیابت

پیش دیابت زمانی رخ می دهد که سطح قند خون فرد بالاتر از حد طبیعی است، اما برای تشخیص دیابت به اندازه کافی نیست. حدود ۷۰٪ افرادی که پیش دیابت دارند استعداد ابتلا به دیابت نوع ۲ و حمله قلبی را نیز دارا هستند. مطالعات زیادی نشان داده است که کاهش وزن ، فعالیت فیزیکی و همچنین اصلاح سبک زندگی در بعضی موارد سطح قند خون را به وضعیت معمولی باز گردانده است [۱۰].

۴. دیابت بارداری^۸

اگر بدن نتواند انسولین کافی برای پاسخگویی به نیازهای اضافی بدن در زمان بارداری را تولید کند این نوع از دیابت رخ می دهد. ممکن است در طول سه ماهه اول بارداری وجود داشته باشد. این نوع دیابت در ۲ تا ۵٪ بارداری ها رخ می دهد و معمولاً بعد از دوران بارداری رفع می شود[۷]. مطالعات نشان داده ۴۰٪ از زنانی که دیابت بارداری را تجربه کرده اند در آینده مبتلا به دیابت نوع ۲ شده اند و همچنین فرزند متولد شده مستعد چاقی و ابتلا به دیابت است[۹، ۱۰].

⁷ Beta Cells

⁸ Gestational Diabetes

۲.۲. علائم، تشخیص و درمان

ابتلا به بیماری دیابت منجر به بروز علائمی در بدن فرد بیمار می شود که در ادامه به تعدادی از آنها اشاره کرده و در مورد تشخیص و درمان دیابت موادردی را مطرح می کنیم.

۱.۰.۲. علائم بیماری دیابت

علائم متداول بروز دیابت شامل، تکرر ادرار، تشنجی زیاد، احساس گرسنگی و اشتهاي بالا، عفونت زخم، ترمیم طولانی مدت زخم، اضطراب، خستگی، کاهش وزن، بیحالی شدید و تاری دید می باشد. برخی از عوامل خطرابتا به دیابت مانند سابقه خانوادگی، اضافه وزن ، چاقی ، فشارخون بالا ، اعتیاد به مصرف دخانیات و مصرف غذایی چرب و شیرین می باشد. دیابت می تواند در ۱ تا ۲٪ موضع از عمل جراحی ، بیماری ژنتیکی خاص، دارو، سوء تغذیه و یا عفونت ناشی شود [۷، ۹].

۲.۰.۲. تشخیص و درمان دیابت

عوامل متعددی مثل معاینه پزشکی، وجود علائم، سابقه پزشکی، آزمایش خون وادرار براساس مقدار قند خون و پروتئین موجود درادرار به تشخیص دیابت کمک می کند. از جمله عواملی که در درمان دیابت موثر است کنترل رژیم غذایی مناسب، انجام فعالیت فیزیکی ، ورزش و همچنین تزریق انسولین (در ۴۰٪ از مبتلایان به دیابت نوع ۲) می باشد. مطالعات نشان داده است که فعالیت فیزیکی به طور چشمگیری می تواند ریسک ابتلا به دیابت نوع ۲ را کاهش دهد [۷، ۹].

dataacademy.ir

۳.۰.۲ آمار ها

بیماری های مسری از جمله دیابت اصلی ترین علت مرگ و میر در جهان و به ویژه در کشور های آسیای شرقی محسوب می شوند. قاره آسیا یک سوم جمعیت جهان را در بر دارد. بیشتر از ۶۰٪ جمعیت دیابتی جهان در کشورهای آسیای شرقی زندگی می کنند . همچنین سبک زندگی و نوع تغذیه مردم در سرتاسر دنیا با هم متفاوت است و همین موضوع علت مشاهده آمار های گوناگون مربوط به دیابت در سرتاسر جهان است. طبق گزارشات^۹ IDF تعداد بیماران دیابتی در سال ۲۰۱۴ برابر ۳۸۷ میلیون نفر بوده که حدود ۲۰۰ میلیون نفر آن مربوط به کشور های آسیای شرقی بوده است. این در حالی است که نیمی از جمعیت دیابتی جنوب شرق آسیا هنوز مورد تشخیص قرار گرفته نشده اند. پیش بینی شده است تا سال ۲۰۴۰ حدود ۶۴۲ میلیون نفر در سراسر جهان مبتلا به بیماری دیابت باشند که این مقدار بیشتر از دو برابر جمعیت مبتلا به دیابت در سال ۲۰۰۸ است. سن ابتلا در حال کاهش است. تعداد زنان دیابتی بیشتر از تعداد مردان دیابتی است و اغلب افراد دارای اضافه وزن می باشند. تا آخر سال ۲۰۱۳ دیابت باعث مرگ ۱.۵ میلیون نفر در جهان شده و همچنین هزینه های پرداخت شده مربوط به این بیماری ۵۴۸ میلیون دلار بوده است [۱, ۸, ۶]. کشور آمریکا در رابطه با افزایش نرخ شیوع دیابت و هزینه های مصرفی آن رشد بسیار بالایی در دهه اخیر داشته است. در سال ۲۰۱۰ دیابت به عنوان هفتمین علت مرگ و میر در آمریکا ثبت شده است. در سال ۲۰۱۱ ۲۰۱۱ دیابت اصلی ترین دلیل از کار افتادگی کلیه در ۴۴٪ از موارد رخ داده شده آمریکا بوده است. دیابت به رگ های ریز خونی در شبکیه آسیب می رساند که باعث تاری دید می شود؛ بین سالهای ۲۰۰۵ تا ۲۰۰۸ حدود ۵.۲۸٪ از افراد ۴۰ سال

^۹ International Diabetes Federation

به بالا، مبتلا به تاری دید چشم ناشی از دیابت بوده اند. در سال ۲۰۱۰ تعداد ۷۳۰۰۰ مورد قطع عضو در افراد ۲۰ سال به بالای مبتلا به دیابت انجام شده است. هزینه های انجام شده مربوط به دیابت اعم از مستقیم و غیر مستقیم در سال ۲۰۱۰ برابر ۲۴۵ میلیون دلار بوده است [۱۰].

۳. داده کاوی پزشکی ، درخت تصمیم

۱۰. داده کاوی پزشکی

فرایند انتخاب ، کاوش و مدلسازی حجم انبوهی از داده جهت کشف الگوهای ناشناخته یا روابطی که نتایج شفاف و مفیدی را فراهم می کند را داده کاوی گویند [۴]. Jianchao Han در مقاله خود در سال ۲۰۰۹ داده کاوی را اینگونه تعریف کرده است: " به کل فرایند بکارگیری متدولوژی کامپیوتری برای گسترش دانش از طریق داده ها ، داده کاوی می گویند "[۱۲]. در دنیای امروز داده های مربوط به بیماران و مریضی های مختلف بسیار گسترده اند، تحلیل و ملاحظات لازم روی تمامی فاکتور های آنها بوسیله شخص پزشک غیر ممکن است. بنابراین به سیستم های هوشمندی برای ملاحظه فاکتورهای مختلف و شناسایی یک مدل مناسب بین پارامتر های مختلف مشاهده شده نیاز است. از طرفی پزشکان و محققان باید قادر به آزمایش و جستجوی اطلاعاتی باشند که به آنها اجازه می دهد تشخیص دقیق تری از بیماری ها داشته باشند. توسعه سیستم های پشتیبانی تصمیم گیری تشخیص پزشکی می تواند به جامعه پزشکی در روند تشخیص و پیش بینی بیماری ها کمک کند. استفاده از سیستم های هوشمند برای تشخیص و درمان بیماری ها، باعث کاهش جدی خطای سیستم پزشکی و در نتیجه کاهش هزینه ها و از دست رفتن جان انسان ها می شود [۱۲، ۱۱].

داده کاوی در حوزه سلامت می تواند جهت فراهم کردن تحلیل داده های مراکز پزشکی برای ارائه منابع بهتر ، شناسایی و پیش بینی سریعتر از بیماری ها استفاده شود، که صرفه جویی در هزینه های ناخواسته و گران آزمایشات مربوط به حوزه سلامت و پزشکی را در بر دارد [۱۴]. داده کاوی در حوزه پزشکی تفاوت هایی با حوزه های غیر پزشکی دارد و با چالش هایی رو به رو است برای مثال کیفیت داده های پزشکی به خاطر مقادیر مفقود ^{۱۰} شده ناشی از بیماریهای مشابه و مریض های مشابه پایین است. سیستم اطلاعات بیمارستان برای اهداف اقتصادی طراحی شده است نه تحقیقاتی. همچنین محقق باید در مورد داده های بیمار، حفظ حریم خصوصی و محرومگی آن طبق ضوابط تعیین شده مطمئن شود. اشتباه در ورود داده در مرحله پیش پردازش و ارزیابی منجر به کشف الگوهای مشکوک و در نتیجه شکایات مردم می شود [۱۵]. در حوزه پزشکی تکنیکهای مختلف داده کاوی اعم از طبقه بندی ، خوش بندی و قوانین انجمنی در مطالعات سال های اخیر استفاده شده اند.

تکنیک های طبقه بندی از دو گام پیروی می کنند. در گام اول (یادگیری یا آموزش ^{۱۱}) : یک مدل طبقه بندی متشکل از قوانین طبقه بندی، با تحلیل داده های یادگیری حاوی برچسب های کلاس ایجاد می شود. در گام دوم (آزمایش ^{۱۲}) : طبقه بند (بوسیله داده های آزمایش) برای میزان صحت و تواناییش جهت طبقه بندی رکوردهای ناشناخته جهت پیش بینی ، آزمایش

¹⁰ Missing Values

¹¹ Train

¹² Test

می شود. مرحله آزمایش در مقایسه با یادگیری بسیار ساده ، با پیچیدگی پائین و نیازمند محاسبات کمتر است. فرایند یادگیری در طبقه بندی به صورت نظارتی و درخوشه بندی به صورت غیر نظارتی انجام می شود. در روش نظارتی برچسب های کلاس پیش بینی بر خلاف روش های غیر نظارتی ، از قبل مشخص است . طبقه بندی جهت دسته بندی داده در برچسب های کلاس از پیش تعریف شده استفاده می شود . کلاس یک خصیصه در مجموعه داده است که در آمار به عنوان یک متغیر وابسته در نظر گرفته می شود. در حوزه پژوهشی تکنیک های طبقه بندی می توانند برای کمک به تشخیص و پیش بینی بیماری های متفاوت از جمله دیابت استفاده شوند. الگوریتم های طبقه بندی مختلفی نظیر طبقه بند بیزین ساده ، شبکه عصبی ، SVM^{۱۳}، رگرسیون، درخت تصمیم و طبقه بند دسته جمعی برای تشخیص و پیش بینی انواع بیماری ها در مطالعات متعددی مورد بررسی قرار گرفته اند. در بین تکنیک ها و الگوریتم های مختلف داده کاوی در این مقاله قصد داریم روی طبقه بند درخت تصمیم تمرکز کنیم.

۲.۰۳. طبقه بند درخت تصمیم

درخت تصمیم یک تکنیک طبقه بندی برپایه یادگیری نظارتی و فرمی جهت بیان یک نگاشت است[8]. توانایی درخت تصمیم در مدلسازی روابط غیرخطی می باشد[۱۶]. در ساختار درخت تصمیم، یک مدل پیش بینی از مجموعه قوانین استخراج شده طبقه بندی بدست می آید[13]. این پیش بینی می تواند در حوزه های مختلف از جمله پژوهشی صورت گیرد.

۱۰.۰۳. یادگیری

dataacademy.ir

درخت تصمیم با شکاف مجموعه داده منبع به زیر مجموعه هایی بر اساس یک آزمایش مقدار ویژگی توسط عمل بازگشتی تقسیم و غلبه برای انتخاب ویژگی ها بواسیله یک تعداد پویش پارامتری روی مجموعه داده، از بالا تا پایین درخت یادگیری می شود. این فرایند روی هر زیر مجموعه مشتق شده به صورت بازگشتی تکرار می شود. عمل بازگشتی وقتی کامل می شود که یا شکاف امکان پذیر نیست یا یک طبقه بند یکتاوی بتواند روی هر عنصر از زیر مجموعه مشتق شده به کارگرفته شود[۱۳، ۱۲].

۲.۰۴. القای درخت تصمیم

درخت شامل نودهای ریشه، میانی(فرزنده) و برگ است و متد انتخاب ویژگی کلید فرآیند ساخت درخت است. خصیصه های انتخاب شده نود های درخت تصمیم هستند. در خت تصمیم فرض می کند که نمونه ها به کلاس های مختلفی تعلق دارند که مقادیر مختلفی حداقل در یکی از خصیصه ها یشان وجود دارد. القای درخت گام یادگیری طبقه بندی است. بهترین ویژگی که داده های یادگیری را تقسیم می کند، بعنوان ریشه درخت در نظر گرفته می شود. نمونه ها با شروع از نود ریشه بر اساس مقادیر ویژگیشان مرتب شده و طبقه بندی می شوند. هر نود در درخت تصمیم نمایانگر یک ویژگی در یک نمونه جهت طبقه بندی شدن است. و بر اساس آن ویژگی و تمام پاسخ های ممکن ، یک سوال دارد . در رابطه با تصمیم گیری برای یک مورد از ریشه شروع می کند و از یک مسیر بر اساس سوالات و پاسخ های نودهای میانی تا زمانی که بالاخره به یک برگ برسد

^{۱۳} Support Vector Machine

پیروی می کند و برچسب برگ برای کلاس آن مورد استفاده می شود. در نهایت هر برگ درخت با استفاده از یکی از کلاس ها برچسب می خورد که نمایانگر نتیجه پیش بینی می باشد [۱۵، ۱۳، ۴]. مسئله القاء طرح ریزی و ایجاد کردن درخت های تصمیم بهینه دو دویی یک مسئله NP کامل است و تئوریسین ها راه های اکتشافی موثر را برای نزدیک شدن به ساخت درختان تصمیم بهینه جستجو می کنند [۳].

3.2.3. الگوریتم های شکاف داده و انتخاب ویژگی

برای پیدا کردن بهترین ویژگی و شکاف داده های مجموعه داده متدهای متفاوتی وجود دارد [۳]. هدف الگوریتم شکاف دهنده پیدا کردن متغیر و آستانه ای است که ترتیب همگن نتایج دو زیر گروه یا بیشتر از نمونه ها را بیشینه می کند. این قدم در نود برگ تا کامل شدن ساختار درخت تکرار می شود [۱۷]. بهترین ویژگی برای شکاف توسط متدهای انتخاب ویژگی انتخاب می شود که بهترین جدا سازی رکورد ها را برای هر یک از برچسب های کلاس انجام می دهد. [۱۵] معمول ترین الگوریتم های ریاضی شکاف مجموعه داده به ۳ دسته تقسیم می شوند که به ترتیب شرح داده می شوند.

1.3.2.2. شکاف داده بر پایه Information Gain

مفهوم Information Gain توسط آقای Hunt و همکارانش در سال ۱۹۶۶ ارائه شد. آقای Quinlan در سال ۱۹۷۹ درخت تصمیم ID3 را ارائه داد. این درخت تصمیم برای انتخاب ویژگی جهت شکاف داده ها و القای درخت تصمیم از استفاده می کند. الگوریتم ID3 یک رویکرد حریصانه است که در آن درخت تصمیم Information Gain برپایه Entropy است. Entropy یک مفهوم معرفی شده در آن درخت تصمیم است که در آن درخت تصمیم گیری به صورت بازگشتی از بالا به پایین بروش تقسیم و غلبه ساخته می شود. درخت با مجموعه یادگیری از تاپل ها و برچسب های کلاس مربوط به آنها ساخته می شود. مجموعه یادگیری با فرایند ساخت درخت تصمیم به صورت بازگشتی به زیر مجموعه های کوچکتر تقسیم می شود. از عمل Entropy برای تقسیم نمونه ها به دو زیر کلاس استفاده شده و میزان همگنی و یکدست بودن مجموعه داده محاسبه می شود؛ اگر مجموعه داده کاملا همگن بود "۰" Entropy می شود. در غیر اینصورت به صورت مساوی تقسیم می شود، تا Entropy مقدار "۱" داشته باشد. Information Gain بستگی به کاهش Entropy دارد. یک ویژگی با بیشترین Entropy مقدار Information Gain را بر می گرداند. بنابراین یک ویژگی با بیشترین Information Gain به عنوان ویژگی شکاف در درخت تصمیم به ریشه منتقل می شود. سپس نودی که بیشترین Information Gain را دارد، والد نسل بعدی می شود. این فرایند تکرار می شود تا زمانی که به نود برگ برسد و درخت تصمیم کامل شود. درخت تصمیم بر اساس پیدا کردن بیشترین شاخه های همگن ساخته می شود. رابطه (۱) نحوه محاسبه Information Gain و رابطه (۲) نحوه محاسبه Entropy را نشان می دهد.

$$\text{Info Gain}(\text{Parent}, \text{Child}) = \text{Entropy}(\text{parent}) - [\sum_{i=1}^m p(c_i) * \log_2 p(c_i)] \quad (1)$$

$$* \text{Entropy}(c_1) + p(c_2) * \text{Entropy}(c_2) \dots]$$

$$\text{Entropy} = \sum_{i=1}^m p(c_i) * \log_2 p(c_i) \quad (2)$$

در رابطه (۱) Child و Parent نمایانگر والد و فرزند می باشند . بر طبق رابطه (۲) احتمال نود فرزند i برابر با p_i است و m تعداد کلاس ها برای ویژگی هدف است[۱۴, ۱۵]. بعد از درخت ID3 آقای Quinlan در سال ۱۹۹۳ الگوریتم C4.5 را ارائه داد که توسعه ای از الگوریتم ID3 است. درخت C4.5 از روش Gain Ratio جهت القا استفاده می کند. همانطور که از رابطه (۳) مشخص است، Gain ratio از تقسیم Splitting Information بر روی Information Gain بدست می آید. برای کاهش تاثیرات تبعیضانه ای که ممکن است به خاطر تعداد زیادی از مقادیر برای یک ویژگی رخ دهد از این الگوریتم شکاف استفاده می شود . یک درخت تصمیم بر اساس Gain Ratio از لحاظ اندازه گیری صحت و مدیریت مسائل بزرگ ، بهتر از Information Gain می باشد. درخت تصمیم با در نظر گرفتن تعداد و اندازه شاخه ها برای یک ویژگی مورد نظر ساخته می شود. Gain Ratio تاثیرات تبعیضانه Information Gain را کاهش می دهد و به صورت یکنواخت مقادیر وسیعی را به یک ویژگی خاص اعطا می کند. یک ویژگی با بیشترین مقدار Gain Ratio بعنوان ویژگی شکاف دهنده برای ساخت درخت تصمیم انتخاب می شود[۱۴]. برای بدست آوردن Splitting Information از رابطه (۴) استفاده می کنیم ؛ که در آن مجموعه داده D به تعداد v پارتیشن مطابق با v خروجی از آزمایش روی صفت خاصه A تجزیه می شود[۱۵]. درخت های تصمیم دیگری مانند C5.0, j48, Gini Ratio نیز از استفاده می کنند.

$$\text{Gain Ratio} = \frac{\text{Information Gain}}{\text{Splitting Info}} \quad (3)$$

$$\text{Split Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} * \log_2 \left[\frac{|D_j|}{|D|} \right] \quad (4)$$

dataacademy.ir

۲.۳.۲.۲. شکاف داده بر پایه Gini Index

Gini Index در سال ۱۹۸۴ توسط آقای Breiman و همکارانش ارائه شد . درخت Gini Index از CART استفاده می کند. این درخت برای مسائل طبقه بندی و رگرسیون استفاده می شود. Gini Index سطح ناخالصی داده های دریافت شده را اندازه گیری کرده ، یک شکاف باینری برای هر متغیر در نظر گرفته و یک درخت دودویی می سازد؛ که خروجی هر نود میانی درخت دقیقا دو کلاس برای ویژگی دریافت شده است. Gini Index برای هر ویژگی محاسبه می شود و سپس بصورت بازگشتی ویژگی با کمترین مقدار بعنوان ویژگی شکاف دهنده داده انتخاب شده و درخت ساخته می شود[۱۶]. برای فهم ساده، درخت تصمیم به مجموعه ای از قوانین IF-Then از ریشه تا برگ تجزیه می شود. قسمت IF بعنوان پیش گویی کننده و قسمت Then بعنوان برچسب کلاس نود برگ (نتیجه مورد انتظار) در نظر گرفته می شود[۱۶]. اگر احتمال A مین کلاس i باشد، برای k کلاس هدف از ویژگی مورد نظر Gini Index به صورت رابطه (۵) محاسبه می شود.

$$\text{Gini Index} = 1 - \sum_{i=1}^k p_i^2 \quad (5)$$

۳.۳.۲.۳. شکاف داده بر پایه Chi-Squared

Chi-Squared test در سال ۱۹۰۰ توسط آقای Pearson ارائه شد[۱۷]. درخت تصمیم CHAID از تکنیک Chi-Squared استفاده می کند. در این درخت از معیار انتخاب صفت بر اساس آزمایش χ^2 استفاده می شود. به طور معمول در

زمینه بازاریابی مستقیم برای انتخاب گروه های مشتریان و پیش بینی چگونگی پاسخ آنها به بعضی از متغیرهایی که تاثیر روی بقیه متغیرها می گذارند، استفاده می شود. CHAID تعامل بین متغیرها را در مجموعه داده با شناسایی گروه های گستته ای از پاسخدهندگان کشف کرده، و در پی کشف و پیش بینی تاثیرات برمتغیرهای وابسته با درنظر گرفتن پاسخدهیشان روی متغیرهای توضیحی می باشد. از آنجاییکه CHAID به داده آماری نیاز دارد، نیازی به گستته سازی متغیرهای عددی ندارد[۱۲]. درخت CHAID در مطالعات مربوط به تشخیص و پیش بینی بیماری توسط داده کاوی بسیار کمتر از بقیه درخت های تصمیم مورد استفاده قرار گرفته است.

۴.۲.۳. بیش برآش^{۱۴} در درخت تصمیم

بک درخت تصمیم یا هر فرضیه یادگیری از h_1 بیش برآش داده های یادگیری گفته می شود اگر یک فرضیه دیگری از h_2 وجود داشته باشد که دارای خطای بزرگتری از h_1 است وقتی روی داده های یادگیری، آزمایش شده باشد؛ و اما دارای یک خطای کوچکتر از h_1 باشد وقتی روی کل مجموعه داده آزمایش شده است. این اتفاق زمانی رخ می دهد که میانگین Entropy کمتر از آستانه تعیین شده است، یا تعداد داده در یک شاخه درخت کمتر از تعداد داده تعیین شده در یک دایرکتوری است[۱۳]. برای جلو گیری از بیش برآش ۲ روش مطرح می شود که به اختصار شرح داده می شود.

- توقف الگوریتم یادگیری قبل از اینکه به نقطه کاملاً برآزند (با اندازه بودن) از داده های یادگیری برسد.
- هرس کردن کامل درخت تصمیم. هرس مشکل بیش برآش مجموعه داده را با متدهای آماری حل می کند. اگر شاخه های مهم بصورت مناسب هرس شوند فهم ساختار داده dataacademy کلی برای افراد ساده می شود.

اگر درخت در هر دو روش دارای صحت پیش بینی یکسان بود، درخت با تعداد برگ کمتر ترجیح داده می شود. دو نوع هرس وجود دارد. پیش هرس در هر مرحله القا درخت تصمیم همراه با فاز ساختن درخت انجام می شود. پس هرس شاخه ها را از یک درخت کامل رشد کرده حذف می کند. برای هر نود غیر برگ در درخت الگوریتم هرس نرخ خطای مورد انتظاری تخمین می زند که اگر زیر درخت در آن نود هرس شد رخ می دهد. بنابراین نرخ خطای مورد انتظار زمانی رخ می دهد که نود هرس نشده با استفاده از نرخ های خطای هر شاخه ای که ترکیب شده بوسیله وزن دهی مطابق نسبت مشاهدات در امتداد هر شاخه تخمین زده شود. اگر هرس کردن نود به نرخ خطای بیشتری از آنچه انتظار میرفت منجر شود، زیردرخت نگه داشته می شود در غیر اینصورت هرس می شود[۱۴].

۵.۰.۳. مزایا و معایب درختان تصمیم

مزایا:

- یکی از خصیصه های مفید درخت های تصمیم قابلیت درک، تفسیر، تجسم مرکز کلاس داده و فهم ساده آن توسط اشخاصی است که آن را مشاهده می کنند. همچنین افراد به راحتی متوجه می شوند که چرا یک نمونه در یک کلاس خاص طبقه بندی شده است و می توانند با این دید که کدام خصیصه تاثیر بیشتری روی یک کلاس داده گذاشته است، کل ساختار داده را تحلیل کنند[۱۵،۳].

^{۱۴} Over Fitting

- درخت تصمیم برای ویژگی‌های گستته جهت ساخت مدل طبقه‌بندی بسیار مناسب است و سرعت قابل توجهی در انجام کار دارد.
- پیاده‌سازی درخت تصمیم ساده است، برای حجم بالایی از داده به خوبی کار می‌کند و نتایج با دقت و صحیحی از طبقه‌بندی را به دست می‌آورد [۱۳].

معایب:

- اغلب درختان تصمیم اگر مسئله نیاز به پارتبیشن بندی قطری یا آریب^{۱۵} داشته باشد خوب کار نمی‌کنند.
- درخت تصمیم می‌تواند برای بعضی از مفاهیم به خاطر وجود تکرار^{۱۶} در یک مسئله و همچنین تعداد ویژگی‌های زیاد در مجموعه داده به شدت نمایش پیچیده ای داشته باشد. برای رفع این مشکل باید عمل هرس انجام شود یا استفاده از الگوریتمی که ویژگی‌های پیچیده در یک نod را با دوری از تکرار و کپی کردن پیاده‌سازی کند [۱۳].

۴. طبقه‌بندی در حوزه سلامت و ارزیابی طبقه‌بندی‌ها

۴.۱. الگوریتم‌های طبقه‌بندی در حوزه سلامت

الگوریتم‌های طبقه‌بندی زیادی در حوزه سلامت وجود دارند که با توجه به معایب و محدودیت‌هایی که دارا هستند ما را دراستفاده درخت تصمیم برای پیش‌بینی بیماری دیابت^{۱۰} می‌نمود. در حوزه داده کاوی پزشکی از شبکه عصبی نیز استفاده گذارد. شبکه عصبی از پس داده‌های نویز دار به خوبی بر می‌آید ولی فرایند یادگیری آن بسیار طولانی است؛ (حدود $\frac{31}{10}$ کند تر از درخت CART). در شبکه عصبی فقدان قدرت توصیف و شفاف سازی مدل برخلاف درخت تصمیم، بدلیل تعداد زیاد نود‌های وزن دهی شده را شاهد هستیم. وزن دهی شبکه بسیار مهم و حساس است و روی نتیجه طبقه‌بندی تاثیر می‌گذارد. از آنجاییکه بصورت جعبه سیاه عمل می‌کند نتایج شبکه عصبی برخلاف درختان تصمیم بوسیله پزشک قابل تفسیر نیست. طبقه‌بند بیزین ساده الگوریتم طبقه‌بندی ساده ای است و سرعت بالایی دارد. یک مشکل آن فرضیه بنیادیش مبنی مستقل بودن خصیصه‌ها از همدیگر است. مثلا در فیلد پزشکی بسیاری از بیماری‌ها و شرایط سلامتی به شدت به هم وابستگی دارند. مثل فشار خون و BMI^{۱۷} که در تشخیص بیماری دیابت به هم وابستگی دارند و در نتیجه طبقه‌بندی می‌تواند خطایجاد کند. در طبقه‌بند SVM نیاز به انتخاب یک تابع کرنل مناسب داریم که خود یک چالش محسوب می‌شود؛ زیرا این طبقه‌بند توابع کرنل زیادی فراهم می‌کند. SVM برای مسائل طبقه‌بندی دو کلاسه طراحی شده، مرحله یادگیری آن بسیار کند است و منابع محاسباتی گستره‌ای نیاز دارد [۱۵].

¹⁵ Diagonal Partitioning

¹⁶ Replication

¹⁷ Body Mass Index

۲۰۴. ارزیابی طبقه بند ها

قبل از ساخت و ارزیابی مدل پیش بینی باید عمل پیش پردازش داده که شامل تمیز کردن داده، حذف مقادیر نویز دار و نرمالیزه کردن داده است انجام شود. پس از آن عمل K-Fold Cross Validation با اندازه های حدوداً یکسان تقسیم می شود. تعداد زیر مجموعه داده تصادفی به K زیر مجموعه ویژه متقابل (از D_k تا D_1) با اندازه های حدوداً یکسان تقسیم می شود. هر دفعه مدل طبقه بندی مجموعه ها معمولاً ۱۰ در نظر گرفته می شود. مدل طبقه بندی K دفعه آموزش و آزمایش می شود. هر دفعه مدل طبقه بندی روی همه داده ها به جز یک قسمت آموزش می بیند و سپس روی تک قسمت باقیمانده آزمایش می شود. این فرایند در شکل (۱) نشان داده شده است. این عمل باعث می شود در مدل پیش بینی بیش برآش رخ ندهد [۸, ۷].



[17] 10-fold Cross validation

بعد از ساخته شدن مدل پیش بینی و آموزش و آزمایش آن، ماتریس اغتشاش^{۱۸} تولید می شود. مقادیر TN, TP, FN, FP به ترتیب نمایانگر False Positive (به اشتباه بیمار پیش بینی شده)، False Negative (به اشتباه سالم پیش بینی شده)، True Positive (بدرسنی بیمار پیش بینی شده) و True Negative (بدرسنی سالم پیش بینی شده) می باشند. مدلی که کمترین مقدار FN را داشته باشد مناسب تر از بقیه طبقه بند ها است و تاثیر زیادی در نتیجه ارزیابی دارد. ماتریس اغتشاش در جدول (۱) نمایش داده شده است.

جدول (۱): ماتریس اغتشاش

		کلاس های پیش بینی شده	
		Yes	No
کلاس های واقعی	Yes	TP	FN
	No	FP	TN

برای ارزیابی و محاسبه کارایی مدل پیش بینی معمولاً از ۳ معیار ارزیابی Accuracy (صحت)، Sensitivity (حساسیت) و Specificity (ویژگی) استفاده می شود که نحوه محاسبه آنها به ترتیب در رابطه های (۶), (۷) و (۸) آمده است.

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (6)$$

¹⁸ Confusion Matrix

$$\text{Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (7)$$

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}} \quad (8)$$

۵. تحلیل مقایسه ای مدل های پیش بینی

در این قسمت از مقاله ۷ کار انجام شده اخیر با تکنیک های مختلف درخت تصمیم در جهت تشخیص و پیش بینی بیماری دیابت با یکدیگر مقایسه می گردد.

۱.۵. مطالعه اول

در سال ۲۰۰۹ آقای Jianchao. H و همکارانش یک مدل پیش بینی بر پایه طبقه بند درخت تصمیم جهت پیش بینی بیماری دیابت ساختند. از مجموعه داده PIMA که شامل اطلاعات افراد دیابتی و غیر دیابتی است استفاده شد. مجموعه داده شامل ۷۶۸ نمونه با ۸ متغیر است که تمام نمونه ها زنان هندی با حد اقل سن ۲۱ سال هستند. ویژگی های استفاده شده در این مطالعه برای پیش بینی بیماری دیابت برابر با تعداد دفعات بارداری بیمار، مقدار آزمایش قند خون ۲ ساعته، فشار خون انبساطی^{۱۹}، ضخامت پوست عضله بازو، میزان انسولین ۲ ساعته، BMI، سابقه خانوادگی دیابت، سن و کلاس پیش بینی است فاز پیش پردازش داده به تنها ۶۰٪ زمان کل پروژه را صرف می کند و تلاش گسترده ای در کل فرایند داده کاوی دارد. بعد از فاز پیش پردازش داده و مدیریت داده های مفقودی ویژگی های ضخامت پوست عضله بازو، میزان انسولین ۲ ساعته از لیست متغیر ها حذف شدند و بقیه متغیر ها گستته سازی شده و به دسته های مختلفی تقسیم بندی شدند. با این کار تمرکز روی زیر مجموعه ها به جای تک تک متغیر ها افزایش پیدا می کند و پیچیدگی تحلیل مدل پیش بینی را بدون از دست دادن صحت کاهش می دهد. در نهایت ۷۲۳ رکورد برای ساخت مدل پیش بینی انتخاب شد. در این تحقیق برای اینکه مقیاس تاثیر هر متغیر روی نتایج بخاطر تفاوت در رنج آن ها استاندارد سازی شود از نرمالسازی-Min-Max استفاده شد؛ برای مثال BMI از رنج ۱۰.۶۷ تا ۲۰.۱۸ و قند خون از ۴۴ تا ۱۹۹ است. درخت تصمیم ID3 انتخاب و به وسیله ابزار Rapid miner مدل طبقه بندی ساخته شد. تحلیل داده توسط نرم افزار Rapid miner نشان داد که بیمار با قند خون بالا بیشترین احتمال برای ابتلاء دیابت را دارد. با این دلیل مقدار قند خون ناشتا^{۲۰} به عنوان نود ریشه درخت در نظر گرفته شد. در مجموعه داده PIMA ۲۴۸ بیمار دیابتی و ۴۷۵ نفر بیمار نبودند. درخت تصمیم ۲۳۱ نفر بیمار دیابتی و ۴۹۲ نفر غیر دیابتی را پیش بینی کرد. از ۲۳۱ نفر ۳۳ نفر اشتباه پیش بینی شد که پیش بینی صحیح به ۱۹۸ نفر کاهش یافت. پس از بدست آوردن ماتریس اغتشاش مقادیر Accuracy و Sensitivity و Specificity به ترتیب برابر ۸۳.۸۰٪ و ۹۳.۰۵٪ بدست آمد.

¹⁹ Diastolic Blood Pressure

²⁰ Fasting Blood Sugar

۲.۵. مطالعه دوم

در سال ۲۰۱۱ آقای Asma. A و همکارانش از طبقه بند درخت تصمیم برای پیش‌بینی دیابت نوع ۲ استفاده کردند. از مجموعه داده PIMA استفاده شده و همه ۷۶۸ نفر از بیماران انتخاب شده منحصرًا زنان هندي با حداقل ۲۱ سال سن هستند. مجموعه داده شامل ۹ ویژگی تعداد دفعات بارداری بیمار، مقدار تست قند خون ۲ ساعته، فشار خون انساطی، ضخامت پوست عضله بازو، میزان انسولین ۲ ساعته، BMI، سابقه خانوادگی دیابت، سن و کلاس پیش‌بینی است. پس از پیش‌پردازش داده و مدیریت داده‌های مفقودی تعداد ۷۲۴ نمونه با ۶ ویژگی قند خون ۲ ساعته، فشار خون انساطی، BMI، سابقه خانوادگی دیابت و سن باقی ماندند.

گسسته سازی متغیرهای عددی جهت کاهش پیچیدگی مسئله انجام شد. برای مثال تعداد دفعات بارداری ۰ یا ۱ برابر کم، بین ۲ تا ۵ برابر متوسط و بیشتر از ۶ بار، زیاد در نظر گرفته شد. برای القای درخت تصمیم از نرم افزار WEKA استفاده شد. مجموعه داده با انتخاب الگوریتم J48 که یک یادگیرنده و پیاده‌سازی از درخت تصمیم C4.5 می‌باشد طبقه بندی شد. در نهایت از ویژگی‌های انتخاب شده و ۷۲۴ نمونه مجموعه داده ۵۶۶ نمونه یا ۱۸.۷۸٪ بدروستی طبقه بندی شد که Accuracy مدل محاسبه شود و ۱۵۸ نمونه یا ۸۲٪ از کل نمونه‌ها اشتباه طبقه بندی شد. Specificity و Sensitivity مدل پیش‌بینی به ترتیب برابر ۸۰.۵۰٪ و ۷۲.۳۳٪ بدست آمد.

dataacademy.ir

۳.۵. مطالعه سوم

در سال ۲۰۱۳ Xue-Hui Meng و همکارانش طبقه بندی‌های رگرسیون، شبکه عصبی و درخت تصمیم را برای پیش‌بینی دیابت و پیش‌دیابت بوسیله فاکتورهای ریسک معمول در این بیماری به کار گرفتند. نمونه‌های مورد مطالعه از مردم کشور چین می‌باشند. مجموعاً ۱۴۸۷ فرد ۲۰ سال به بالا در این مطالعه شرکت داده شده که اطلاعات آنها از سال ۲۰۰۷ تا ۲۰۰۸ جمع‌آوری شده است. ۷۳۵ نفر دیابتی (با قند خون ناشتا بالا) و پیش‌دیابتی و ۷۵۲ نمونه معمولی (بوسیله چک آپ فیزیکی در ۲ سال گذشته) تأیید شدند. مصاحبه استاندارد برای بدست آوردن تاریخچه خانوادگی دیابت، اندازه گیری‌های مربوط به بدن و فاکتورهای ریسک سیک زندگی انجام شد. سپس مدل‌های پیش‌بینی با ۱۲ متغیر ورودی و ۱ متغیر خروجی بوسیله اطلاعات استاندارد دریافت شده از مصاحبه ایجاد شد. رضایت‌کننده‌ها ها قبل از جمع‌آوری داده بدست آمد. داده‌های یادگیری شامل ۷۰٪ شرکت کنندگان (۱۰۳۱ مورد) و داده‌های آزمایش ۳۰٪ شرکت کنندگان (۴۵۶ مورد) را تشکیل دادند. مدل پیش‌بینی با ابزار SPSS Modeler ver14.1 طرح ریزی شد. ویژگی‌های مهم در مورد نمونه‌ها در این تحقیق از جمله سن، جنسیت، وضعیت تاہل، سطح سواد، تاریخچه خانوادگی دیابت، مصرف قهوه، میزان خواب، BMI، فعالیت فیزیکی، استرس کاری، مصرف غذاهای شور و مصرف ماهی می‌باشند. مقدار BMI با وزن به کیلوگرم تقسیم بر مربع قد به متر (m^2) محاسبه شد؛ و BMI بیشتر مساوی ۲۵ به عنوان اضافه وزن در نظر گرفته شد. متغیر مستقل (خروجی) یک دسته بندی با اینری دارد (۰ یعنی معمولی و ۱ یعنی دیابتی). در این مطالعه فاکتورهای ریسک سبک زندگی شامل متغیرهای مختلفی می‌باشد که استرس کاری یکی از آنها محاسبه شود؛ بنابراین کل تعداد متغیرهای

ورودی ۱۹ مورد شد. از درخت ۰ C5. بر پایه Information Gain جهت ساخت مدل پیش بینی استفاده شده و درخت تصمیم بهترین نتیجه را در بین دو طبقه بند دیگر با Accuracy، Sensitivity و Specificity بترتیب ۷۷.۸۷٪، ۷۷.۸۰٪ و ۷۵.۱۳٪ بدست آورد. سن بیشتر، تاریخچه خانوادگی دیابت، BMI بالا، مصرف غذای شور، فعالیت بدنی کم و میزان استرس مربوط به کار با این بیماری ارتباط مثبت داشت در حالی که سطح تحصیل و نوشیدن قهوه ارتباط منفی با این بیماری داشت. بقیه ویژگی ها از جمله مصرف سیگار، مصرف الکل، خوردن سبزیجات و میوه ارتباط مستقیم با این بیماری داشتند.

۴.۵. مطالعه چهارم

در سال ۲۰۱۴ آقای b . Saba و همکارانش درخت های تصمیم C4. 5 , Id3, CART را به صورت یک مجموعه طبقه بند ترکیبی برای پیش بینی بیماری دیابت بوسیله دو مجموعه داده Uci,Biostate به کار گرفته اند. طبقه بندی های Majority Voting, Adaboost, Bayesian Boosting, Bagging and Stacking استفاده شده در این مطالعه از تکیک های جمعی استفاده شده است که کمک به شناسایی هدف اصلی این مطالعه شناسایی بهترین چارچوب طبقه بندی دسته جمعی برای درختان تصمیم است که کمک به شناسایی موثر بیماران دیابتی و مهتمراز آن صحت بالای بدست آمده دارد. درخت های تصمیم مختلف با معیارهای شکاف مختلف از جمله Information Gain, Gain Ratio , Gini Index بر عنوان طبقه بندی های پایه استفاده شدند. هر طبقه بند بر اساس داده های آموزش مرحله یادگیری را اجرا می کند. ایده اساسی طبقه بندی های دسته جمعی، وزن دهی چندین طبقه بند مجزا و ترکیب آنها برای بدست آوردن نتیجه ایست که بهتر و شایسته تر از هر یک از طبقه بندی ها می باشد . در این مطالعه از ابزار Rapid Miner5 استفاده شد. بوسیله Fold cross validation ۱۰ از ۹۰٪ داده ها بعنوان یادگیری و ۱۰٪ داده ها برای آزمایش استفاده شد. نتایج بدست آمده نشان داد که بهترین مجموعه طبقه بند ترکیبی Bagging می باشد. برای ارزیابی کلی طبقه بند مقادیر Accuracy و Sensitivity و Specificity به ترتیب برابر با ۹۱.۵۶٪، ۹۵.۶۳٪ و ۹۳.۳٪ و ۶۸.۶٪ بدست آمد.

۵.۵. مطالعه پنجم

در سال ۲۰۱۴ توسط O. Purnik و همکارانش برای شناسایی افرادی که ریسک پایینی در ابتلا به دیابت دارند بوسیله پایگاه داده TLGS که مربوط به مرکز قندخون و چربی شهر تهران است یه مدل پیش بینی کننده ساختند. اطلاعات افراد طی ۴ فاز در سالهای مختلف جمع آوری شد: فاز ۱: ۱۹۹۹-۲۰۰۱، فاز ۲: ۲۰۰۲-۲۰۰۵، فاز ۳: ۲۰۰۵-۲۰۰۸ و فاز ۴: ۲۰۰۹-۲۰۱۲. تعداد ۱۵۰۰۰ نفر با سن بیشتر مساوی ۳ سال از ۱۳ منطقه مختلف شهر تهران با متند نمونه برداری خوش تصادفی انتخاب شدند. افرادی که از قبل دیابت داشتند حذف شد و ۱۰۳۱۰ نفر بدون دیابت باقی ماندند. اطلاعات افرادی که دارای مقدیر مفقودی درمورد مقدار قند ناشتا و قند دو ساعته بودند از مجموعه داده حذف شدند. در نهایت ۶۶۴۷ نفر غیر دیابتی برای ساختن مدل پیش بینی انتخاب شدند. بوسیله پرسشنامه و مصاحبه اطلاعاتی از قبیل: سن ، جنسیت ، وضعیت تاہل ، میزان تحصیلات، مصرف دخانیات ، فعالیت های بدنی تاریخچه پزشکی و دارویی جمع اوری شد. اندازه قدم، وزن، دور کمر، ران، مچ و BMI بر طبق استاندارد بدست آمد. فشار خون از بازوی راست با فشار سنج جیوه ای بدست آمد و قند خون ناشتا بوسیله آزمایش خون اندازه گیری شد. سپس نمونه های خون جهت تحلیل در آزمایشگاه مرکز تحقیقات ۳۰ تا ۴۵ دقیقه

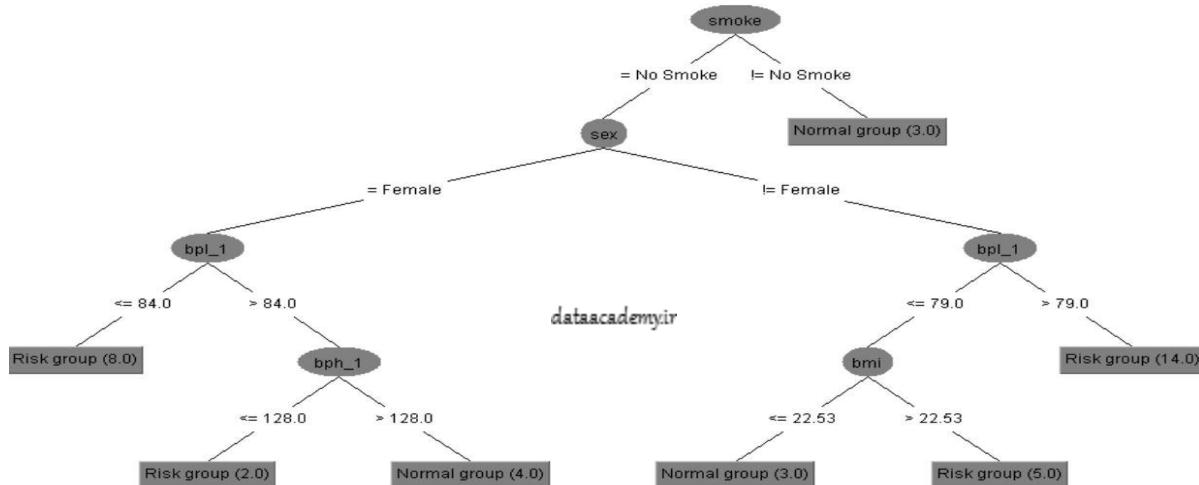
در سانتریفیوژ نگهداری و خنک شد . برای کلی سازی، داده سطح ابتدایی داده با مفاهیم سطح بالا جایگزین شد؛ برای مثال تحصیلات به ۳ سطح کمتر از ۵ سال ، ۱۲-۵ سال و بیشتر از ۱۲ سال تقسیم بندی شد. مجموعه داده نهایی به ۶۰ متغیر ورودی (پیش بینی کننده ها) و یک متغیر نتیجه که به ۲ دسته (رخداد دیابت نوع ۲ و عدم رخداد دیابت نوع ۲) تعریف شده است، تقسیم شد. برای ساخت مدل پیش بینی از ابزار Konstanz Information Miner استفاده شد. برای معيار انتخاب ویژگی استفاده یادگیری ۳۰٪ از نمونه ها و برای آزمایش انتخاب شدند. از درخت CART و Gini Index برای معيار انتخاب ویژگی استفاده شد. درخت تصمیم با مجموعه یادگیری شامل ۴۶۵۲ رکورد ساخته شد و با ۱۹۹۵ رکورد، آزمایش شد. برای اولین متغیر جدا سازی، سطح قند خون ناشتا انتخاب شد. نتایج ارزیابی مدل پیش بینی برابر با ۹۰.۵۰٪ Accuracy و ۳۱٪ sensitivity داده شد. در دیابت نوع ۲ روی افراد غیر دیابتی است. درخت تصمیم در این مطالعه روی درختان تصمیم جهت پیش بینی احتمالات رخ ندادن دیابت نوع ۹۷٪ و ۹٪ specificity داشت. تنها ۷۲۹ نفر (۱۱٪) بیماری دیابت ، شناسایی شد.

۶.۵ مطالعه ششم

در سال ۲۰۱۵ آقای Nongyao.n و همکارانش از طبقه بند های رگرسیون، درخت تصمیم، بیزین ساده و شبکه عصبی جهت طبقه بندی ریسک بیماری دیابت استفاده کردند و همچنین از تکنیک های Bagging and Boosting برای توانمندی طبقه بندها استفاده شد. هدف این مطالعه پیش بینی بیماری دیابت برای هر فرد بدون نیاز به آزمایش خون و رفتن به بیمارستان بوده است. مجموعه داده اولیه از ۲۶ واحد مراقبت اصلی (PCU) در بیمارستان SPR تایلند در بین سالهای ۲۰۱۲-۲۰۱۳ جمع آوری شده است. هر فرد فرم گزینشی را پر می کند که برای شناسایی گروه ریسک دیابت در این مطالعه استفاده شده است. متغیر های ورودی خروجی به وسیله مشاوره با شخص پزشک تعیین می شوند و سپس از اطلاعات عمومی در فرم گزینش انتخاب می شوند. در این مطالعه ۱۱ متغیر ورودی و یک متغیر خروجی دو بخشی در نظر گرفته شده است. متغیر های ورودی به ترتیب جنسیت، BMI، قد، دورکمر، وزن، فشارخون انقباضی^{۲۱}، فشار خون انبساطی ، سابقه فشار خون در خانواده ، سابقه دیابت در خانواده، مصرف الکل، سیگار کشیدن و متغیر خروجی متغیر قند خون ناشتا است که تقسیم می شود به دو گروه معمولی و ریسکی. گروه معمولی افرادی هستند که مقدار قند خون ناشتا آنها کمتر از ۱۰۰ میلی گرم در دسی لیتر است و گروه ریسکی بین ۱۰۰ تا ۱۲۵ می باشد. در این تحقیق افراد بالای ۱۲۵ لحظه نشده اند چرا که آنها در گروه دیابت طبقه بندی شده و تحت درمان می باشند. سپس از حذف رکوردهای دارای مقادیر مفقودی از مجموعه داده ، مجموعه داده نهایی شامل ۳۰۱۲۲ نفر شد که به دو گروه معمولی ۱۹۱۴۵ نفر و ریسکی ۱۰۹۷۷ نفر تقسیم شد. با ترکیب الگوریتم های Bagging, Boosting, Random Forest و ۴ طبقه بند ذکر شده ۱۳ مدل طبقه بندی برای پیش بینی ساخته شد. الگوریتم Bagging با تولید تصادفی یک مجموعه داده یادگیری شروع می کند. سپس جایی که یک کلاس پیش بینی ایجاد شده است را مدل می کند. این فرایند چندین بار تکرار می شود تا هر دفعه عمل جایه جایی داده ها صورت گیرد. پیش بینی خروجی نهایی به وسیله رای اکثریت مدل های پیش بینی شده ارائه می شود. الگوریتم Boosting با به کارگیری یک وزن به تمام مشاهدات مجموعه داده در حال یادگیری شروع می شود. سپس طبقه بندش (مثلا درخت تصمیم) مدلسازی می شود.

²¹ Systolic Blood Pressure

سپس این گام ها چندین بار تکرار می شود و کلاس های پیش بینی از رای اکثریت ترکیبات می باشد. الگوریتم Random Forest از الگوریتم Bagging مدل سازی شده است. این روش بصورت بالقوه ای می تواند صحت طبقه بندی را بهبود بخشد و همچنین روی مجموعه داده هایی با تعداد بسیار زیادی از متغیر های ورودی بخوبی کار می کند. الگوریتم با ساختن یک ترکیبی از درختانی که هر کدام برای یک کلاس رای داده اند شروع می شود. در این مطالعه درخت تصمیم بر پایه Information Gain ساخته شد. سپس مدل با 10-Fold cross validation آزمایش شد. بیشترین Accuracy به ترتیب برای مدل های Random Forest و درخت تصمیم با مقدار ۸۵.۵۵٪ و ۸۵.۰۹٪ و کمترین مقدار مربوط به طبقه بند ترکیبی Bagging و بیزین ساده با ۸۰.۹۶٪ بوده است. در شکل(۲) قسمتی از درخت تصمیم ساخته شده در این مطالعه نمایش داده شده است. عدد داخل پرانتز در برگ های درخت احتمال تعلق گرفتن نمونه به کلاس پیش بینی مورد نظر است



شکل(۲): مثالی از درخت تصمیم [8]

۷.۵. مطالعه هفتم

در سال ۲۰۱۵ آقای Mahmoud^{۲۲} و همکارانش از طبقه بندی های درخت تصمیم، بیزین ساده، SVM و Knn^{۲۳} شبکه عصبی برای تشخیص بیماری دیابت نوع ۲ در شهر تبریز استفاده کردند. هدف این تحقیق کمینه کردن هزینه اقتصادی سربار از خطای طبقه بندی دیابت بوسیله تکنیک ها و الگوریتم های مختلف در سیستم سلامت، با در نظر گرفتن هزینه های غلط طبقه بندی افراد سالم به عنوان بیمار و بالعکس بوسیله پارامتر های نرخ TP و FP بوده است. مجموعه داده دارای ۲۵۳۶ رکورد که برای این بیماری آزمایش شده اند، می باشد. داده این بیماران بوسیله دانشگاه علوم پزشکی تبریز طی ۳ ماه اول سال ۲۰۱۰ جمع آوری شد. تعداد دیابتی ها ۲۳۰۵ نفر و غیر دیابتی ها ۲۳۱ است. بنابراین مجموعه داده نامتعادل است. BMI متغیر یا ویژگی در مجموعه داده با مشورت پزشکان و متخصصان انتخاب شد که عبارتند از جنسیت، سن، وزن، قد، سابقه خانوادگی دیابت، سابقه بارداری، سابقه دیابت بارداری، سابقه سقط جنین، سابقه فشار خون بالا، سابقه مصرف دارو برای

²² K-Nearest Neighbor

فشار خون بالا، میزان فشار خون انبساطی و میزان فشار خون انقباضی. از Information Gain برای ساخت درخت تصمیم و برای پیاده سازی و آزمایش تکنیک های طبقه بندی از نرم افزار WEKA ورژن ۳.۶ استفاده شد. برای تحلیل کارایی هر طبقه بند از 10-fold cross validation استفاده شد. در کاربردهای پزشکی میانگین هزینه طبقه بندی اشتباه FN بیشتر از FP است و مدل طبقه بندی مناسب تر است که کمترین مقدار FN را داشته باشد. زیرا پیش بینی منفی اشتباه برای فرد بیمار عواقب خطروناکی را برای اوی در پی دارد. از آنجاییکه مقدار Sensitivity مهم تر از Specificity یا FPR^{23} است ما نیاز به کمینه کردن مقدار FPR به کمترین مقدار قابل قبول برای TPR^{24} داریم. ۵nn و درخت تصمیم به ترتیب بیشترین Specificity را با مقدار ۹۷.۸۵٪ و ۹۹.۸۷٪ داشتند. شبکه عصبی و درخت تصمیم کمترین FP را با Accuracy به ترتیب به ۶۶.۹۷٪ و ۹۶.۸٪ دارا بودند و ۵nn و درخت تصمیم دارای کمترین FN بودند. بیشترین Accuracy به ترتیب به شبکه عصبی و درخت تصمیم با ۹۷.۴۴٪ و ۹۵.۰۳٪ اختصاص یافت.

8.5 جدول مقایسه

در جدول (۲) مقایسه تحلیلی ۷ مقاله بررسی شده از نظر تکنیک شکاف داده، نوع درخت تصمیم، تعداد نمونه و ویژگی ها و همچنین ارزیابی آنها نمایش داده است.

جدول (۲): مقایسه تحلیلی طبقه بندی های درخت تصمیم در مقالات مطرح شده

مقاله	سال	داده	شکاف	تکنیک	درخت تصمیم	تعداد نمونه	متغیر	آزار داده کاوی	Specificity	Sensitivity	Accuracy
۱	۲۰۰۹		ID3	Information Gain		723	Input:6 Output:1	Rapid Miner	93.05%	79.83%	88.50%
۲	2011		C4.5:J48	Gain Ratio		724	Input:6 Output:1	WEKA	80.50%	72.33%	78.18%
۳	2013		C5.0	Information Gain		1487	Input:19 Output:1	Spss modeler ver14.1	75.13%	80.68%	77.87%
۴	2014		Information Gain Ratio, Gini Index	Ensemble Bagging: CART C4.5 , ID3		-	-	Rapid Miner ver. 5	68.33%	95.63%	91.56%
۵	2014		Gini Index	CART		6647	Input:60 Output:1	KNIME Ver.2.6	97.90%	31.10%	90.50%
۶	2015		Information Gain			۳۰۱۲۲	Input:11 Output:1	-	-	-	85.09%
۷	2015		Information Gain			۲۵۳۶	Input:13 Output:1	WEKA ver3.6	66.97%	97.08%	95.03%

²³ False Positive Rate

²⁴ True Positive Rate

۶. نتیجه گیری

در این مقاله ابتدا به تعاریف مربوط بیماری دیابت، انواع آن و آمارها اخیر پرداخته شد. سپس در مورد داده کاوی در حوزه پزشکی، انواع درختان تصمیم، روشاهای القای آنها و معایب و مزایای آنها توضیحاتی داده شد. به مقایسه ۷ مورد از کارهای اخیر در زمینه پیش‌بینی دیابت بوسیله طبقه بند درخت تصمیم پرداخته شد. در مطالعه اول تعداد رکوردها و ویژگی‌های کمی بررسی شده است ولی با خاطر پیش‌پردازش مناسبی که روی داده‌ها انجام شده، صحت مناسبی از مدل را شاهد بودیم. در مطالعه دوم تعداد نمونه‌ها یکی بیشتر از مطالعه اول بوده و تعداد ویژگی‌ها برابر مطالعه اول بوده است ولی با تغییر تکنیک شکاف داده حدود ۱۰٪ کاهش صحت طبقه بند را شاهد بودیم. یکی از دلایل آن می‌تواند خصوصیت Gain Ratio باشد که در مجموعه داده‌های کوچک صحت کمتری نسبت به مجموعه داده‌ها و مسائل بزرگ‌تر دارد. در مطالعه سوم به همراه افزایش تعداد نمونه‌ها نسبت به دو مطالعه قبلی تعداد فاکتورهای ریسک نیز افزایش پیدا کردند. بعضی از این فاکتورها مثل سطح تحصیل و نوشیدن قهوه ارتباطی با بیماری دیابت نداشتند. یکی از دلایل کاهش صحت طبقه بندی نسبت به ۲ کار قبلی می‌تواند افزایش تعداد فاکتورهای نامرتب به بیماری در مدل پیش‌گیری باشد. در مطالعه چهارم از طبقه بند مجموعه ترکیب استفاده شده است که یک طبقه بند کار آمد به حساب می‌آید و صحت خوبی را بدست آورده است. در این مطالعه از تعداد نمونه‌ها و ویژگی‌ها صحتی نشده است و از این جنبه نمی‌توانیم اظهار نظری داشته باشیم. در مطالعه پنجم نمونه‌های مجموعه داده افراد غیر دیابتی بوده اند و هدف پیش‌بینی دیابت برای افرادی بوده است که ریسک پایینی در ابتلا به این بیماری دارند. بنابراین تعداد افراد غیر دیابتی بالا باعث افزایش مقدار TN و کاهش مقدار TP می‌شود. مقدار پایین FN تاثیر مستقیمی روی کاهش Sensitivity dataacademy.ir می‌گذارد و میزان آن با مقدار ۱۱.۱٪ ثبت می‌شود. افزایش تعداد متغیرها نسبت به مقالات قبلی تأثیری در افزایش صحت نداشته است ولی در مورد افزایش تعداد نمونه‌ها نمی‌توانیم این ادعا را داشته باشیم. در مطالعه ششم ۱۳ طبقه بند مختلف مورد آزمایش قرار گرفت که درخت تصمیم با اختلاف بسیار ناچیز دومین صحت بالای این مطالعه را بدست آورد. از مطالعه ششم نتیجه می‌گیریم که طبقه بندی در ترکیبی تأثیر مثبتی در بهبود نتایج طبقه بندی و پیش‌گیری الگوریتم‌های طبقه بندی مجزا دارند. مطالعه هفتم با اختلاف ۰.۲٪ و میزان ۹۵.۰۳٪ دومین صحت بالای طبقه بندی را در مقایسه با شبکه عصبی در مطالعه مربوطه بدست آورد. این عدد بیشترین مقدار در مقایسه تحلیلی ما در این مقاله می‌باشد. تفاوت ویژگی‌ها و فاکتورهای ریسک در مقاله هفتم با بقیه کارهای بررسی شده استفاده از ویژگی‌هایی مربوط به دیابت بارداری می‌باشد. در نهایت نتیجه می‌گیریم که افزایش تعداد ویژگی‌های مورد بررسی تأثیر زیادی در افزایش صحت نداشته، افزایش تعداد نمونه‌ها می‌مورد بررسی در کل در بهبود صحت پیش‌گیری موثر است و بهترین روش شکاف داده با توجه به ارزیابی کلی، Specificity Sensitivity، Accuracy، Information Gain

مراجع:

1. D. A. G. Idf. (2015), "Update of mortality attributable to diabetes for the IDF Diabetes Atlas: Estimates for the year 2013," *Diabetes research and clinical practice*, vol. 109, p. 461.
2. M H. Nadimi ,M.Taki and F.Habibollahi.(2014), *Data Mining Concepts And Applications* vol. 1.
3. T. N. Phy, "Survey of classification techniques in data mining.(2009)," in *Proceedings of the International MultiConference of Engineers and Computer Scientists*, pp. 18-20.
4. X. -H. Meng, Y. -X. Huang, D. -P. Rao, Q. Zhang, and Q. Liu.(2013),"Comparison of 3 data mining models for predicting diabetes or prediabetes by risk factors," *The Kaohsiung journal of medical sciences*, vol. 29, pp. 93-99.
5. R. M. Goodman and P. Smyth.(1990),"Decision tree design using information theory," *Knowledge Acquisition*, vol. 2, pp. 1-19.
6. F. Aguiree, A. Brown, N. H. Cho, G. Dahlquist, S. Dodd, T. Dunning, et al .(2013),"IDF diabetes atlas,".
7. R. W. Grant and M. S. Kirkman.(2015),"Trends in the Evidence Level for the American Diabetes Association's "Standards of Medical Care in Diabetes" From 2005 to 2014," *Diabetes care*, vol. 38, pp. 6-8.
8. N. Nai-arun and R. Moungmai.(2015),"Comparison of Classifiers for the Risk of Diabetes Prediction," *Procedia Computer Science*, vol. 69, pp. 132-14.
9. C. f. D. Control, Prevention, C. f. D. Control, and Prevention.(2011),"National diabetes fact sheet on diabetes and prediabetes in the United States," *Atlanta, GA: US, Centers for Disease Control and Prevention*, vol. 201.
10. A. D. Association, "National Diabetes Statistics Report.(2014)," *Estimates of diabetes and its burden in the epidemiologic estimation methods*. *Natl Diabetes Stat Rep*, pp. 2009-2012.
11. A. A. Al Jarullah, (2011)"Decision tree discovery for the diagnosis of type II diabetes. " (IIT), International Conference on , pp. 303-307.
12. J. Han, J. C. Rodriguez, and M. Beheshti, "Discovering decision tree based diabetes prediction model.(2009)" in *Advances in Software Engineering*, ed: Springer, pp. 99-109.
13. M. Heydari, M. Teimouri, Z. Heshmati, and S. M. Alavinia,(2015)"Comparison of various classification algorithms in the diagnosis of type 2 diabetes in Iran," *International Journal of Diabetes in Developing Countries*, pp. 1-7.
14. S. Bashir, U. Qamar, F. H. Khan, and M. Y. Javed, "An Efficient Rule-Based Classification of Diabetes Using ID3, C4. 5, & CART Ensembles.(2014)," (FIT),*12th International Conference on*, pp. 226-231.
15. I. Yoo, P. Alafaireet, M. Marinov, K. Pena-Hernandez, R. Gopidi, J. -F. Chang, et al . , "Data mining in healthcare and biomedicine: a survey of the literature.(2012)" *Journal of medical systems*, vol. 36, pp. 2431-2448.
16. A. Ramezankhani, O. Pournik, J. Shahrabi, D. Khalili, F. Azizi, and F. Hadaegh, "Applying decision tree for identification of a low risk population for type 2 diabetes. Tehran Lipid and Glucose Study.(2014)" *Diabetes research and clinical practice*, vol. 105, pp. 391-398.
17. D. Delen, G. Walker, and A. Kadam,(2005)"Predicting breast cancer survivability: a comparison of three data mining methods," *Artificial intelligence in medicine*, vol. 34, pp. 113-127.
18. R. L. Plackett, "Karl Pearson and the chi-squared test,(1983)" *International Statistical Review/Revue Internationale de Statistique*, pp. 59-72.