

داده کاوی بر پایه روش‌های شبکه عصبی و درخت تصمیم در تشخیص زود هنگام ریسک ابتلا به دیابت بارداری

مریم میرشریف^{۱*}، سعید روحانی^۲

• پذیرش مقاله: ۹۶/۳/۳۰

• دریافت مقاله: ۹۶/۲/۱۰

مقدمه: امروزه در دنیای مدرن صنعتی خطر ابتلا به بیماری‌های مزمن به طرز چشمگیری افزایش یافته است. دیابت بارداری یکی از مسائل مهم در حوزه سلامت است و در صورتی که درمان نشود مشکلات و عوارض جانبی متعددی برای مادر و فرزندش به همراه دارد. این پژوهش به دنبال پیش‌بینی ریسک و هشدار به موقع در ابتلا به دیابت بارداری به مادر می‌باشد تا در اوایل بارداری از ابتلا جلوگیری به عمل آید.

روش: این پژوهش که به صورت کاربردی-پیمایشی انجام شد و از دو رویکرد شبکه عصبی و درخت تصمیم در داده‌کاوی به منظور تجزیه و تحلیل آزمایشی داده‌ها و پیش‌بینی استفاده گردید. داده‌های استخراج شده نرمال‌سازی شده و پس از آماده‌سازی در نرم‌افزار Matlab تجزیه و تحلیل شدند.

نتایج: تحقیق حاضر در پی یافتن پاسخ به این پرسش است که "آیا دو روش داده‌کاوی شبکه عصبی و درخت تصمیم در تشخیص به هنگام و درست ریسک ابتلا به دیابت بارداری از صحت لازم برخوردار است؟" و می‌توان از آن‌ها برای تشخیص درست استفاده نمود؟ نتایج تحقیق نشان می‌دهد که روش‌های داده‌مدار در بهبود صحت و درستی پیش‌بینی مؤثرند، در کشف دانش ضمنی و تشخیص روابط پنهان بین داده‌ها عملکرد مناسبی دارند و خطای تصمیم‌گیری در هر دو روش در حد قابل پذیرش و بسیار به هم نزدیک است. **نتیجه‌گیری:** نتایج تحقیق حاکی از آن است که از رویکردهای داده‌مدار می‌توان در مراکز درمانی و سایر بیماری‌های کمتر شناخته شده استفاده نمود و پیشگیری به موقع، مدیریت خود بیمار و کاهش هزینه‌های درمانی را میسر ساخت.

کلیدواژه‌ها: داده‌کاوی، شبکه‌های عصبی هوشمند، درخت تصمیم، دیابت بارداری، تشخیص

ارجاع: میرشریف مریم، روحانی سعید. داده‌کاوی بر پایه روش‌های شبکه عصبی و درخت تصمیم در تشخیص زود هنگام ریسک ابتلا به دیابت بارداری. مجله انفورماتیک سلامت و زیست پزشکی ۱۳۹۶؛ ۴(۱): ۶۸-۵۹.

۱. کارشناس ارشد مدیریت فناوری اطلاعات، دانشگاه علوم و تحقیقات تهران، تهران، ایران.

۲. دکترای مهندسی سیستم، استادیار، دانشکده مدیریت، دانشگاه تهران، تهران، ایران.

* نویسنده مسئول: تهران، دانشگاه علوم و تحقیقات تهران

مقدمه

دیابت بارداری قندی نوعی از دیابت است که برای نخستین بار در بارداری ظاهر می‌شود [۱]. در مناطق مختلف جهان شیوع متفاوتی دارد که محدوده ۴ تا ۱۷ درصدی برای آن گزارش شده است [۲]. این بیماری یکی از مسائل مهم در حوزه سلامت است در صورتی که درمان نشود مشکلات و عوارض جانبی متعددی برای مادر و فرزندش به همراه دارد. جهت پیشگیری از عوارض ناخواسته‌ای که همواره گریبان‌گیر آنان است نیازمند توجه جدی و هشدار به موقع در ماه‌های اولیه بارداری می‌باشد. به دلیل تغییر در سبک زندگی در شهرهای صنعتی و کاهش تحرک، افزایش چاقی و تعداد افراد دیابتی، فعالیت بدنی کمتر و افزایش سن ازدواج و سایر عوامل محیطی مؤثر، احتمال ابتلا به دیابت نسبت به گذشته بسیار بیشتر شده است [۳]. در سال‌های اخیر سیستم‌های تشخیصی برای بیماری‌هایی همچون دیابت بسیار گسترش یافته که در شرایط مختلف می‌توان از آن‌ها بهره‌مند شد. با توجه به گسترش دسترسی به پایگاه‌های داده مراکز درمانی و داده‌های پزشکی، روش‌های داده‌کاوی همچون درخت تصمیم و شبکه‌های عصبی بسیار مورد توجه قرار گرفته است [۴].

دیابت بارداری به طور معمول در اواسط دوران بارداری (هفته ۲۴-۲۸) تشخیص داده می‌شود که بنا به شدت آن مادر باید تحت مراقبت‌های درمانی همچون تزریق انسولین قرار گیرد [۵، ۶]. در رویکرد پیشنهادی هشدار احتمال ابتلا به دیابت بارداری در ماه‌های اولیه به مادر داده می‌شود تا با توجه به آن و رعایت رژیم غذایی مناسب و مراقبت‌های پزشکی موردنیاز از ابتلا جلوگیری به عمل آید. برای تشخیص به هنگام و در ست امکان ابتلا به دیابت بارداری روش‌های داده‌کاوی شبکه‌های عصبی و درخت تصمیم برای کاهش سطح خطا و افزایش صحت و دقت تشخیص به کار گرفته شده تا پیش‌بینی را بهبود بخشد.

دیابت بارداری

دیابت حاملگی شایع‌ترین عارضه و اختلال متابولیک دوران بارداری است که اغلب بدون علامت است. بیش از نیمی از زنان مبتلا به این بیماری سرانجام دچار بیماری دیابت واضح می‌شوند و پس از ۵ تا ۱۰ سال به دیابت نوع ۱ و یا با شیوع بیشتر به دیابت نوع ۲ مبتلا می‌شوند [۳]. ابتلا به این بیماری صدمات فراوانی را به دنبال دارد و اگر این بیماری پیشگیری و درمان نشود، عوارضی بر روی مادر و جنین دارد که عبارت‌اند از: کاهش

شدید قندخون نوزاد در چند ساعت اول زندگی، بزرگی جثه جنین، صدمات زمان تولد و آسیب‌های مغزی، زردی نوزاد، افزایش مرگ‌ومیر در دوران بارداری و مرگ ناگهانی جنین در ماه‌های آخر در شکم مادر [۴].

برخی عوارض دیابت بارداری بر روی مادران عبارت‌اند از: استرس عصبی، عفونت به طوری که شیوع عفونت‌های ادراری در زنان باردار مبتلا به دیابت بالاتر است، بالا بودن گلوکز مادر که منجر به نارسایی جفت می‌شود، فشارخون بالا، افزایش خطر مسمومیت حاملگی، ورم، افزایش مایع آمنیوتیک، صدمات زایمانی، زایمان قبل از ترم و افزایش عمل سزارین [۷]. در مطالعه‌ای که در بیمارستان طالقانی در مورد پیامد حاملگی در ۵۲ حامله دیابتی انجام شد چنین نتیجه‌گیری شد که میزان سزارین در گروه دیابت آشکار، بیش از سایرین بود که به دلیل ماکروزومی جنین در آن‌ها می‌باشد (۳۲/۵ درصد) [۴]. در چهارمین کنگره بین‌المللی تازه‌های دیابت که در تالار امام بیمارستان امام خمینی (ره) برگزار شد، متخصصین این رشته از کشورهای آمریکا، فرانسه، کانادا و دانمارک و ایران، از افزایش آمار دیابت در کشورهاشان خبر داده و توانمندسازی سیستم سلامت برای خدمت‌رسانی به بیماران دیابتی را با توجه به افزایش دیابت ضروری خواندند و آن را در کنترل این بیماری مؤثر دانستند. دیابت سالانه موجب مرگ بیش از دو میلیون نفر در جهان می‌شود و این آمار، لزوم مداخله جدی برای کنترل این بیماری را می‌طلبد. لیکن پیشگیری از آن آسان‌تر و کم‌هزینه‌تر از درمان آن خواهد بود.

داده‌کاوی

داده‌کاوی یا کشف دانش از میان پایگاه‌های اطلاعاتی علمی است که برای تصمیم‌گیری‌های هوشمندانه با توجه به پیشرفت روز افزون تکنولوژی اطلاعات بسیار کاربرد دارد. تکنیک‌های داده‌کاوی برای پیدا کردن الگوهای جالب در تشخیص پزشکی و درمان مورد استفاده است. انواع مختلفی از تکنیک‌های داده‌کاوی در استخراج اطلاعات که می‌تواند برای پیش‌بینی مورد استفاده قرار گیرد وجود دارد که از آن جمله می‌توان به شبکه‌های عصبی، درخت تصمیم و رگرسیون لجستیک اشاره کرد [۸، ۹].

یکی از روش‌های مناسب برای پیش‌بینی و تشخیص در حوزه‌های پزشکی رویکرد داده‌کاوی می‌باشد. داده‌کاوی روشی داده‌مدار و بر مبنای یادگیری و کشف الگوی پنهان در میان داده‌های حقیقی می‌باشد که از این الگو جهت پیش‌بینی برای موارد مشابه استفاده می‌کند [۱۰].

مطالعه قرار گرفت و با استفاده از شبکه‌های عصبی با الگوریتم بهینه بیزین به کار گرفته شد [۱۸]. در تحقیق مشابهی برای تشخیص دیابت بارداری توسط میرشریف و همکاران به طراحی یک سیستم خبره فازی و رویکرد مدل فازی پرداخته شده، از سیستم خبره استنتاج فازی برای تشخیص دیابت بارداری در شرایط عدم دسترسی به پزشک استفاده شده است. سیستم خبره فازی روشی بر پایه مدل می‌باشد که نیازمند جمع‌آوری دانش افراد خبره در پایگاه قوانین می‌باشد. در مدل سیستم خبره فازی قوانین در پایگاه دانش تحت نظر متخصصان، افراد خبره و کتاب‌های مرجع پزشکی استخراج و جمع‌آوری شده که می‌توان از آن به صورت تصمیم‌یار استفاده نمود. سطح خطای MSE (mean squared error) در سیستم استنتاج فازی حدود ۰/۲۲۷ برآورد گردید [۱۹].

در تحقیق حاضر، پس از تحلیل داده‌های اولیه و داده‌های خروجی بیماران مورد مطالعه، با افرادی مواجه شدیم که در ماه‌های اولیه بارداری علائم هشدار دهنده اولیه یا ریسک فاکتورهای مؤثر در ابتلا به دیابت بارداری همچون قند خون بالا یا وزن و سن بالا در هنگام بارداری نداشته‌اند، با این وجود در انتهای بارداری به دیابت مبتلا شدند. به منظور پوشش‌دهی این دسته از بیماران و تشخیص درست، می‌توان از روش‌های داده‌کاوی بهره جست. روش‌های داده‌کاوی قادرند دانش و روابط پنهان میان حجمی از داده را درک کرده و از این دانش برای پیش‌بینی‌های بعدی استفاده نمایند. روابطی که شاید کمتر کسی متوجه ارتباط معنادار ورودی‌ها با نتایج خروجی شده باشد و نیازمند تحلیل‌های پیچیده‌ای است که ذهن بشر قادر به انجام آن نیست [۱۰]. این تحقیق به منظور کاهش خطا در تشخیص و پیش‌بینی به موقع ابتلا به دیابت در زنان باردار انجام گرفته تا با آنالیز داده‌های کلینیکی با ابزارهای مناسب به کشف دانش و الگوهای روابط بین داده‌ها و در نتیجه تصمیم‌سازی مناسب‌تر و بهبود مراقبت‌های موردنیاز در دوران بارداری منجر شود.

روش

از آنجاکه در این تحقیق به دنبال پیشنهاد روشی برای تشخیص درست و بهنگام دیابت در اوایل بارداری بودیم، از حیث نوع هدف کاربردی می‌باشد و در آن داده‌های کمی از پایگاه‌های داده پزشکی جمع‌آوری شده. در روش‌های درخت تصمیم و شبکه عصبی در داده‌کاوی می‌توان با آزمایش و تکرار به پاسخ بهینه دست‌یافت، از این رو روش تجزیه و تحلیل داده‌ها آزمایشی

تحقیقات بسیاری در زمینه تشخیص با استفاده از شبکه‌های عصبی و سیستم‌های خبره صورت گرفته است. در مقاله‌ای با عنوان مدل‌سازی یک سیستم خبره برای تشخیص دیابت بارداری قندی بر اساس عوامل خطر به مدل‌سازی یک سیستم خبره برای تشخیص دیابت بارداری، با استفاده از معماری شبکه عصبی روبه جلو پرداخته است که ورودی‌های در نظر گرفته شده در آن، با ورودی‌هایی که در این تحقیق به عنوان عوامل مؤثر یا ریسک فاکتورهای مؤثر در پیدایش بیماری در نظر گرفته شده‌اند، متفاوت می‌باشد [۱۰]. سیستم خبره دانش را از فرد خبره و متخصص کسب می‌کند. این دانش را می‌توان با استفاده از متغیرهای کلامی بیان نمود تا فهم و انتقال حالات مختلف در آن آسان‌تر باشد [۱۱]. برای برآورد زمان واقعی و مقدار مناسب تزریق انسولین در دیابت نوع T1DM (type 1 diabetes mellitus)، سیستمی بر اساس مدل کنترل پیش‌بینی غیر خطی NMPC (nonlinear model predictive control)، در سال ۲۰۱۱ توسعه داده شد. این مقاله شامل شبکه عصبی و یک الگوریتم بر اساس منطق فازی می‌باشد که برای کنترل آنالیز و انطباق پارامترهای کنترل NMPC توسعه داده شده است [۱۲]. در مقاله‌ای با عنوان "تشخیص دیابت قندی با استفاده از شبکه‌های عصبی"، تشخیص دیابت با رویکرد شبکه‌های عصبی با الگوریتم پی‌انتشار خطا ارائه شده و شامل هفت ورودی: سن، شاخص توده بدنی، انسولین سرم، گلوکز پلاسما و... می‌باشد [۱۳]. در مقاله دیگری از الگوریتم ژنتیک و سیستم استنتاج فازی عصبی برای تشخیص هوشمندانه دیابت نوع ۱ استفاده شد [۱۴]. سیستم استنتاج فازی عصبی ترکیب دو روش شبکه عصبی و سیستم استنتاج فازی است و شبکه‌ای تطبیق‌پذیر و قابل آموزشی است که به لحاظ عملکرد کاملاً مشابه سیستم استنتاج فازی است و از قابلیت آموزش شبکه عصبی برای وزن‌دهی به پارامترها استفاده می‌کند و در عین حال از قابلیت رابط کاربر مناسب و تعریف مفاهیم فازی در سیستم استنتاج فازی بهره می‌برد [۱۵]. مقالات و تحقیقاتی در تشخیص دیابت با استفاده از شبکه عصبی انجام گرفته که بر مبنای فاکتورهای مؤثر در پیدایش بیماری بنا نهاده شده‌اند [۱۶، ۱۷]. برای تشخیص دیابت قندی بر اساس عوامل خطر با ۱۶ لایه ورودی، مقاله‌ای توسط Sumathy و همکاران منتشر شده که به سهولت تشخیص و کمک به بیماران برای پیش‌بینی ابتلا به دیابت توسط خود فرد کمک می‌کند [۱۷]. همچنین در مقاله Nguyen و همکاران به منظور تشخیص هایپوگلیسمی یا کاهش قندخون در تعدادی کودک مبتلا به دیابت نوع ۱ مورد

دیابتی آمریکا (American Diabetic Association) ADA می‌باشد [۱۰].

داده‌های استخراج شده نرمال سازی شده و پس از آماده سازی در نرم افزار Matlab برای پیاده سازی در شبکه عصبی و در نرم افزار clementine برای ترسیم درخت تصمیم وارد گردید. از ۸۰ درصد داده‌ها برای آموزش و از ۲۰ درصد برای تست استفاده شد و برای سنجش کارایی روش‌های به کارگرفته شده، صحت یا درصد پاسخ‌های درست محاسبه گردید. در ادامه به شرح مختصری از مراحل کار می‌پردازیم.

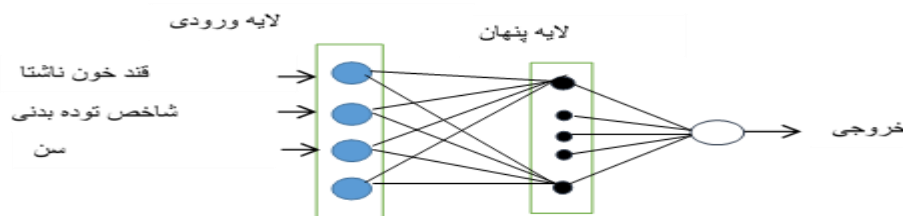
شبکه‌های عصبی مصنوعی

شبکه‌های عصبی سیستم‌های دینامیکی هستند که با پردازش روی داده‌های تجربی، دانش یا قانون نهفته در ورای داده‌ها را به ساختار شبکه منتقل می‌کنند. این سیستم‌ها براساس محاسبات روی داده‌های اولیه و جدید ارائه شده به آن قوانین کلی را فرا می‌گیرند و سیستم‌هایی هوشمندی هستند. پردازش اطلاعات در شبکه‌های عصبی روشی مشابه مغز انسان دارد. این شبکه از تعداد زیادی از عناصر پردازشی به هم پیوسته (نورون‌ها) تشکیل شده که به صورت موازی کار می‌کنند [۱۲]. شبکه‌های عصبی می‌توانند با تحلیل داده‌ها و ارتباط بین ورودی و خروجی‌های واقعی، جنبه‌های پنهان علم و روابط پنهان در پس داده‌ها را کشف نماید. این دانش از وزن‌هایی که به شبکه عصبی اختصاص می‌دهد قابل استنتاج است [۲۰].

شبکه‌های عصبی مصنوعی از معماری سه لایه شامل: لایه ورودی، پنهان و خروجی تشکیل شده است. هر لایه شامل تعدادی نورون یا گره می‌باشد. مدل پرسپترون چند لایه کاربرد موفقیت‌آمیزی در حل برخی مسائل از جمله شناسایی الگو و تخمین تابع دارد [۲۱]. شمایی از شبکه سه لایه به کار برده شده در این مقاله با ورودی و خروجی‌هایش در شکل ۱ نشان داده شده است.

می‌باشد. ابتدا ورودی‌های حقیقی که اطلاعات اولیه بالینی هر بیمار می‌باشد به همراه خروجی حقیقی از ابتلا یا عدم ابتلای شخص به دیابت در پایان بارداری برای انجام آنالیزهای لازم در داده‌کاو فراهم شد. ورودی‌های شبکه عصبی و درخت تصمیم برای داده‌کاو، اطلاعات کد شده پزشکی هستند که از سوابق بیماران استخراج شده‌اند. در هر بیماری‌ای عوامل متعدد و بی‌شماری در شدت و ضعف آن نقش دارند که این عوامل از نظر درجه اهمیت آن در حصول نتیجه با هم متفاوت‌اند [۱۰]، این عوامل در درجه اهمیت و توجه متفاوت‌اند. بر اساس پارامترهای فیزیولوژیکی، پیشینه تحقیقات و استناد بر کتاب‌های مرجع پزشکی و نظر پزشکان متخصص زنان، سه عامل قند خون ناشتا، شاخص توده بدنی (BMI (Body Mass Index)، فشار خون و سن حاملگی از بالاترین درجه اهمیت در میان تمام عوامل مؤثر در بروز دیابت بارداری برخوردارند [۷]. تعداد هزار پرونده بیماران از سال ۱۳۹۰ تا ۱۳۹۳ در تحقیقی میدانی و نمونه‌گیری هدفمند از کلینیک پزشکی تخصصی زنان در تهران مورد بررسی قرار گرفت که از میان آن‌ها دویست و چهل و هشت مورد دارای اطلاعات کامل از ابتدا تا انتهای بارداری بودند (اطلاعات موردنیاز تحقیق میزان فاکتورهای مؤثر قند خون، وزن و سن بارداری در ماه‌های اولیه بارداری و همچنین جواب آزمایش شربت گلوکز در ماه‌های پایانی و هفته بیست و هشتم می‌باشد)، سپس اطلاعات موردنیاز اولیه شامل فاکتورها و عوامل مؤثر بر پیدایش دیابت بارداری (ورودی‌ها: سن و وزن و قندخون) و همچنین نتایج حاصل از تست تحمل گلوکز (Oral Glucose Challenge Test (OGCT) در هفته ۲۴-۲۸ از بارداری (شاخصی برای تعیین خروجی سیستم) از پرونده‌ها استخراج گردید. نتایج تست (OGT (Oral Glucose Test) بر اساس استانداردهای مختلف می‌تواند تحلیل شود.

روش تحلیلی به کار گرفته شده در این مقاله مربوط به انجمن



شکل ۱: معماری شبکه‌های عصبی سه لایه

وزن، فرزند چندم بودن، سابقه فامیلی، عوامل محیطی، فشارخون و غیره، سه عامل قند خون ناشتا در اوایل بارداری و سن بارداری

در هر بیماری عوامل متعددی در پیدایش آن دخیل می‌باشند. در بین فاکتورهای اصلی در پیدایش دیابت بارداری همچون سن،

و وزن با استناد به کتب مرجع پزشکی در بالاترین اولویت برای توجه قرار دارند [۲۲]. به منظور پیش بینی با استفاده از شبکه عصبی مصنوعی که با استفاده از نرم افزار Matlab پیاده سازی گردید، ۳ پارامتر استخراج شده از پرونده بیماران به عنوان ورودی برای آموزش به شبکه اعمال شد. نرون ها در لایه ورودی مشخصات بالینی بیمار می باشد. برای آموزش شبکه، پس از معرفی داده های ورودی، مقادیر خروجی به ازای ورودی ها به سیستم معرفی شده و آموزش شبکه آغاز می شود. پس از آموزش شبکه با استفاده از ورودی و خروجی های دنیای واقعی، می توان از شبکه عصبی برای گرفتن نتیجه در ازای ورودی های جدید بهره برد. هر ورودی دارای یک مقدار حقیقی از خروجی می باشد و در شبکه عصبی هر واحد خروجی یک تابع خطی از ورودی آن است. در شبکه های عصبی مصنوعی توابع خطی پارامترهایی دارند که می توان آن ها را برای رسیدن به حداقل خطا در آموزش تنظیم نمود. می توان شبکه را توسط مجموعه ای از نمونه ها آموزش داد و این کار را تا پیدا کردن بهترین مقدار پارامترها انجام دهیم تا در نهایت به کمترین میزان خطا در شبکه برسیم [۲۳].

برای پیش بینی ریسک ابتلا به دیابت در مرحله اولیه بارداری از یک شبکه روبه جلوی با ناظر با الگوریتم پس انتشار خطا در نرم افزار متلب استفاده شده است. نرون ها می توانند از توابع محرک متفاوتی جهت تولید خروجی استفاده کنند که از رایج ترین آن ها می توان به توابع لگاریتم سیگموئیدی، تانژانت سیگموئیدی و تابع محرک خطی اشاره کرد. تابع تحریک مورد استفاده در این تحقیق از نوع لگاریتم سیگموئیدی (۱) می باشد که مقادیری بین صفر و یک تولید می نماید.

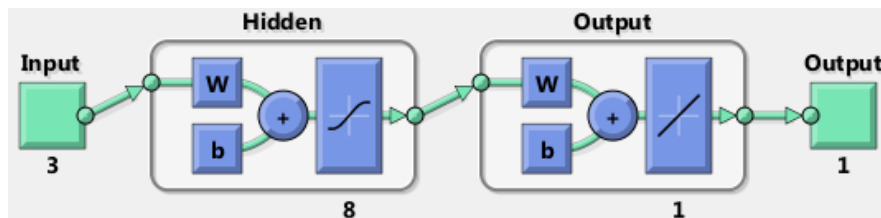
$$f(n) = sig(n) = \frac{1}{1 + e^{-n}} \quad (۱) \text{ فرمول}$$

از میان روش های مختلف آموزش به روش پس انتشار خطا، الگوریتم لونبرگ مارکوارت، به دلیل همگرایی سریع تر در آموزش شبکه های با اندازه متوسط، برای استفاده در تحقیق حاضر انتخاب شده است. در این روش از مشتق اول (گرادیان) و مشتق دوم موسوم به هسین برای اصلاح پارامترها استفاده می شود. الگوریتم پس انتشار خطا، وزن های شبکه و مقادیر بایاس را در جهتی تغییر می دهد که تابع عملکرد با سرعت بیشتری کاهش یابد [۲۴]. مشخصه هدف این مسئله پیش بینی ریسک ابتلا به

دیابت بارداری است که با مقادیر ۰ و ۱ به معنای ریسک کم و ریسک ابتلای خیلی زیاد، مقداردهی شده است؛ بنابراین تنها یک نرون برای لایه خروجی وجود خواهد داشت. با نتیجه خروجی (۱) وجود خطر بالای ابتلا به بیمار اطلاع رسانی می شود. در صورتی که فرد از این هشدار چشم پوشی کند، ممکن است در ماه های باقی مانده از بارداری در معرض ابتلا قرار گیرد. این در حالی است که در صورت توجه به این هشدار و با در نظر گرفتن یک رژیم غذایی مناسب و تحت نظر قرار گرفتن بیمار، می توان از ابتلا و عوارض حاد در پی آن جلوگیری به عمل آورد. شبکه با مقادیر مختلف آموزش داده شده و میزان خطا در پایان هر تکرار اندازه گیری شده است و بهترین پاسخ ها استخراج می شود. تعداد گره در لایه ورودی و خروجی بستگی به مورد تحت بررسی دارد. تعداد نرون ها در لایه میانی در اکثر شبکه های (Multi-Layer MLP (Perceptron با آزمایش و خطا در تکرار الگوریتم آموزش به دست می آید که مناسب ترین روش در اکثر تحقیقات علمی معرفی شده است. n تعداد نرون ها در لایه مخفی است. یکی از مشکلات مهم و دشوار تعیین تعداد بهینه لایه پنهان و تعداد گره ها است. روش آزمون و خطا روشی است که می توان با آن بهترین ساختار شبکه را تعیین نمود. برای پیدا کردن تعداد گره ها در لایه پنهان و داشتن حداقل خطا، الگوریتمی توسط Hirose و همکاران پیشنهاد گردید که با تغییر تعداد گره پنهان به صورت پویا و رسیدن به حداقل مقادیر خطا ادامه می یابد (هر تکراری که در آن $MSE(Mean Squared Error)$ یا متوسط مربع خطا کمتر از مقدار از پیش تعیین شده باشد) [۲۳]. به ازای هر مقداری از n شبکه دارای مقادیر رگرسیون (R) و حداقل مربعات $MSE (Mean Squared Error)$ متفاوتی است که این میزان به ازای $n=8$ بهترین پاسخ را برای سیستم فراهم نمود. همان طور که در جدول ۱ نشان داده شده است، تعداد $n=8$ برای تعداد مناسب گره های پنهان در شبکه پیشنهادی در نظر گرفته شده است و این شبکه دارای کمترین میزان خطای مجذور مربعات و بیشترین مقدار R می باشد. پس از مقایسه مقادیر پیش بینی شده با مقادیر واقعی، مقادیر 0.91 برای R و مقدار 0.07 برای MSE یا مقدار خطا به دست آمده است. که نشانگر عملکرد مطلوب برای پیش بینی با استفاده از شبکه عصبی مصنوعی می باشد. توان دوم همبستگی خطی را (R^2) یا ضریب تعیین می نامند و روابط با ضریب تعیین بالاتر و میانگین مربع خطای کمتر به عنوان مدل برتر انتخاب می شوند [۲۳].

جدول ۱: نتیجه تغییر در تعداد گره پنهان

	MSE[۰,۱]	R ²
n=۸	۰/۰۷۳۸	۰/۹۱۹۸
n=۹	۰/۱۵۴۳	۰/۷۷۴۴
n=۱۰	۰/۱۳۶۵	۰/۵۹۲۱



شکل ۲: بهترین معماری شبکه

درک برای انسان از روابط موجود در یک مجموعه داده‌ای هستند و پیش‌بینی خود را در قالب قوانینی که از نظر پارامترهای آماری برازش مناسبی دارند ارائه می‌کنند. این روش یادگیری برای توابع گسسته و داده‌های خطادار به کار می‌رود و به کشف دانش کمک می‌کند [۲۶]. درخت تصمیم یکی از ابزارهای داده‌کاوی در ارتباط با حل مشکلات بیماری‌ها در دنیای واقعی است. درخت تصمیم مجموعه داده‌های دیابت را طبقه‌بندی می‌کند. یک درخت تصمیم‌گیری یک مدل طبقه‌بندی مناسب با استفاده از مجموعه‌ای از داده‌ها می‌باشد [۲۷] ویژگی‌ها از میان مجموعه‌ای از داده‌ها برای کمک به تصمیم‌گیری استخراج می‌شوند. توزیع کلاس به صورت کلاس با مقدار ۱ به معنی ابتلا به دیابت بارداری و کلاس ۰ به معنی دیابت حاملگی منفی می‌باشد. مراحل الگوریتم درخت تصمیم‌گیری به شرح زیر است: پس از وارد کردن ورودی‌ها از مجموعه داده‌ها، برای هر ویژگی اطلاعات حاصل از تقسیم بر روی آن ویژگی را نمایش داده و بهترین را از میان صفات پیدا می‌نماید. در نتیجه انتخاب گره تصمیم‌گیری و تقسیم بر روی بهترین ویژگی انجام می‌شود. پس از انتخاب نود تصمیم‌گیری یک الگوریتم بازگشتی بر روی زیر شاخه‌ها و تقسیم شدن شاخه‌ها در بهترین ویژگی انجام شده و آن گره به عنوان گره فرزند به درخت اضافه می‌شود [۲۴]. در این مقاله، از الگوریتم درخت تصمیم C4.5 برای داده‌کاوی استفاده شده است. داده‌کاوی هسته مرکزی فرآیند کشف دانش از پایگاه داده KDD (Knowledge Discovery in Database) می‌باشد که به پیدا کردن الگوها و قواعد از بین پایگاه داده‌های بزرگ کمک می‌کند [۲۸]. درخت تصمیم قوانین تصمیم‌گیری و

پس از انتخاب بهترین معماری و ساخت شبکه با ۸ نرون در لایه پنهان (شکل ۲)، فرآیند آموزش و تست مدل آغاز می‌شود. پس از آماده‌سازی داده‌ها مدل شبکه عصبی مصنوعی از نوع شبکه عصبی پیشخور ایجاد گردید. شبکه از نوع شبکه‌های MLP، با قانون یادگیری خطای تصحیح BP (Back Propagation) از نوع آموزش با سرپرست می‌باشد [۲۱]. مدل یادگیری پس انتشار خطا (BP) یک مدل نرم‌افزاری پرکاربرد و موفق است و در میان تمام مدل‌های موجود بیشتر مورد توجه تحقیقات قرار گرفته است [۲۵].

روش جستجو گرادیان است که خطای سیستم و یا نرخ مربعات خطا (MSE) را به حداقل می‌رساند. هر الگویی از ورودی به شبکه ارائه و آموزش شبکه آغاز می‌شود. BP ارزش گره‌های خروجی را تعیین می‌کند و در گره‌های پنهان خطا را تعیین می‌نماید، پس از آن به به‌روزرسانی وزن‌ها بر مبنای کاهش خطای محاسبه شده در هر مرحله می‌پردازد. شبکه‌های عصبی مصنوعی میل طبیعی ذخیره‌سازی دانش تجربی برای استفاده مجدد از آن در آینده دارند. شبکه عصبی می‌تواند داده‌ها و دانش گذشته را ذخیره کند و پس از یادگیری آن را برای استفاده مجدد در دسترس قرار دهد [۲۳].

پس از آموزش شبکه، خطای سیستم از روی خروجی مشخص داده شده توسط کاربر (خروجی واقعی) و خروجی تعیین شده از سوی مدل تعیین می‌گردد [۲۴].

درخت تصمیم‌گیری

یکی از طبقه‌بندی‌های محبوب و متداول برای دسته‌بندی و پیش‌بینی می‌باشد که پیاده‌سازی آن آسان، ساده و تفسیر نتایج آن امکان‌پذیر است. درختان تصمیم قادر به تولید توصیفات قابل

جدول ۳: مقادیر صحت در درخت تصمیم و شبکه عصبی

Accuracy	BIAS	MAD	MSE
۰/۹۰۶	۰/۰۰۰۴	۰/۰۱۶۱	۰/۰۷۴

با توجه به مقادیر خروجی برای میزان خطای محاسبه شده برای شبکه‌های عصبی مصنوعی و درخت تصمیم، نمایانگر درصد خطای قابل پذیرشی برای هر دو روش می‌باشد که بسیار به هم نزدیک هستند. روش‌های کاربردی فراوانی برای پیش‌بینی‌های دقیق مانند شبکه عصبی مصنوعی، سیستم‌های خبره، رگرسیون خطی و لجستیک و درخت تصمیم‌گیری وجود دارد. به نظر می‌رسد که روش‌های داده‌کاوی مورد استفاده در این مقاله (شبکه‌های عصبی و درخت تصمیم‌گیری) دارای مقادیر قابل قبولی از سطح خطا و دقت برآورد در پیش‌بینی برخوردارند و این مقدار در مقایسه با نرخ خطای محاسبه شده در نتایج سایر مقالات برای پیش‌بینی در تشخیص‌های پزشکی، گویای مقادیر قابل‌پذیرش برای کاربرد روش‌های داده‌کاوی را دارد که می‌توانند به عنوان ابزاری جهت بهبود پیش‌بینی استفاده شوند.

بحث و نتیجه‌گیری

مطابق با جدول ۳ که بیانگر صحت دو رویکرد شبکه‌های عصبی و درخت تصمیم است می‌توان چنین نتیجه‌گیری کرد که هر دو روش از مقادیر قابل پذیرش و نزدیک به هم (۹۰ درصد و ۹۳ درصد) برخوردارند و هر دو روش از توانمندی و کارایی لازم در امر تشخیص برخوردارند. پیشنهاد می‌شود در مسائل پیچیده‌ای که امکان دسترسی به داده‌های حقیقی در تحقیق وجود دارد از تکنیک‌های داده‌کاوی و الگوریتم‌های اکتشافی در کشف جنبه‌های پنهان دانش و روابط بین متغیرها بهره‌مند شویم. در مقایسه با تحقیقات مشابه قبلی، موارد موفقیت‌آمیز در تشخیص دیابت قندی و برخی بیماری‌های قلبی و گوارشی یافت می‌شود؛ اما کمتر مقاله‌ای به تشخیص دیابت بارداری توجه نموده است. در بین روش‌های داده‌کاوی، درخت تصمیم از کاربردهای فراوانی برخوردار است. از موارد استفاده از درخت تصمیم در پیش‌بینی می‌توان به «آنالیز مدل درخت تصمیم در دیابت قندی» در سال ۲۰۱۶ اشاره نمود که در آن ادعا شده با استفاده از آنالیز درخت تصمیم و دسته‌بندی j48 که در نرم‌افزار Matlab, weka انجام گرفت، صحت ۹۹/۸۷ درصدی به دست آمده است. از آنجا که دیابت قندی یک بیماری مزمن است بهتر است قبل از ابتلای افراد به آن پیشگیری صورت گیرد و می‌توان با آنالیز ژن‌ها

کنترل را برای داده‌های پر از نویز فراهم می‌کند و می‌توان قوانین و الگوهای منظم در پس پایگاه‌های داده بزرگ را درک نمود، زیرا ساختار درختی شکل می‌تواند به راحتی توسط کاربر نهایی درک شده و به تصمیم‌گیری دقیق‌تر کمک نماید. با تحلیل اطلاعات با نرم افزار Clementine(C4.5) و بررسی اثر سه ورودی بر خروجی، پارامتر سن از ضریب تأثیر بالاتری بر روی خروجی برخوردار بود و ریشه درخت تصمیم بر اساس سن به زیرشاخه‌ها تقسیم گردید. دقت طبقه‌بندی نتایج در مدل درخت تصمیم‌گیری حدود ۰/۹۳۰ برآورد شده که کارایی قابل توجهی را برای مدل نشان می‌دهد.

نتایج

به طور معمول در شبکه‌های عصبی جهت سنجش عملکرد این شبکه‌ها از شاخص‌های ارزیابی ویژه‌ای استفاده می‌شود. نتایج ارزیابی شاخص‌های مذکور به صورت جدول ۲ است.

جدول ۲: نتایج شاخص‌های عملکرد در شبکه عصبی

Method	Accuracy
DT(Clementine)	۰/۹۳۰
ANN(Matlab)	۰/۹۰۶

میزان صحت یا عملکرد شبکه بیانگر تعداد پیش‌بینی‌های درست به تعداد کل پیش‌بینی‌ها می‌باشد و صحت بیشتر از ۸۸ درصد حاکی از عملکرد خوب در پاسخ‌دهی بوده و لذا عملکرد شبکه طراحی شده قابل قبول است. بایاس، مد و کمترین مربع خطا (فرمول‌های ۲، ۳، ۴) شاخص‌های دیگری هستند که بزرگی خطا را نمایش می‌دهند و برای تست اعتبارسنجی مدل می‌توان آن‌ها را از روابط زیر محاسبه نمود. روشی برای پیش‌بینی مناسب‌تر است که مقادیر این شاخص‌ها در آن کمتر باشد.

$$MAD = \frac{1}{n} \sum |actual - forecast| \quad \text{فرمول (۲)}$$

$$BIAS = \frac{1}{n} \sum (actual - forecast) \quad \text{فرمول (۳)}$$

$$MSE(\theta^{\wedge}) = E[(\theta^{\wedge} - \theta)^2] \quad \text{فرمول (۴)}$$

مقادیر محاسبه شده از صحت عملکرد در دو روش شبکه عصبی و درخت تصمیم در جدول ۳ نشان داده شده است که بیانگر عملکرد قابل قبول از دو روش در کاربرد را دارد.

سابقه فامیلی و محل زندگی و فشارخون را با افزایش تعداد نمونه‌ها گسترش داد. در انتها نتیجه‌گیری می‌شود تکنیک‌های داده‌کاوی توانایی تحلیل پایگاه‌های داده بزرگ‌تر و آنالیز روابط بین داده‌ها را دارد و می‌توان در صورت دسترسی به داده‌ها در مراکز درمانی، تحلیل‌ها و نتایج خوبی را استخراج نمود و از آن‌ها به صورت تصمیم‌یار در تشخیص و پیش‌بینی استفاده نمود. پیشنهاد می‌شود از روش‌های داده‌کاوی در مورد سایر بیماری‌هایی که تشخیص زود هنگام در درمان آن حائز اهمیت است استفاده نمود. از آنجا که زنان باردار مبتلا به دیابت بارداری ممکن است به دیابت نوع دوم در زندگی آینده خود مبتلا شوند و آمارها از افزایش سالانه تعداد افراد مستعد به ابتلا گواهی می‌دهند، لزوم توجه و پیشگیری از آن قابل تأمل است. در صورت استفاده فرد از سیستم‌های تشخیصی هوشمند و استفاده از تحلیل‌های مناسب در کشف روابط در بیماری‌های پیچیده، بیمار می‌تواند به بررسی خود بدون مراجعه حضوری با پزشک به خصوص در مناطق دور افتاده پرداخته و از ابتلا به دیابت بارداری و عوارض حاد آن جلوگیری به عمل آورد، که در نتیجه، مدیریت خود بیمار و پیش‌بینی به موقع را میسر ساخته و با هشدار و پیشگیری به موقع به کاهش هزینه‌های پزشکی و درمانی کمک نموده و مانع از افزایش تعداد مادران و کودکان دیابتی در جامعه می‌شود و این مهم خود یکی از تأثیرات و فواید فناوری اطلاعات در عصر مدرن می‌باشد.

و پیشینه دیابت از ابتلا به آن جلوگیری به عمل آورد [۲۹]. همچنین می‌توان به کاربرد موفقیت‌آمیز استفاده از شبکه عصبی در پیش‌بینی و کاربردهای پزشکی همچون مقایسه الگوریتم‌های داده‌کاوی در تشخیص دیابت قندی در سال ۲۰۱۴ اشاره نمود که از روش‌های داده‌کاوی درخت تصمیم و نزدیک‌ترین همسایگی و انفیس (Adaptive Neuro-Fuzzy Inference System) بهره برده که نتایج آن حاکی از صحت بیشتر روش انفیس (فازی-عصبی) در تشخیص و در حدود ۸۰ درصد می‌باشد [۳۰]. همچنین در مقاله‌ای با عنوان "آنالیز تکنیک‌های مختلف داده‌کاوی در تشخیص دیابت قندی" صحت پیش‌بینی در درخت تصمیم حدود ۸۶ درصد و در شبکه عصبی حدود ۷۴ درصد تخمین زده شد [۳۱]. برای داده‌های مختلف با استفاده از درخت تصمیم، شبکه‌های عصبی و استفاده از سایر روش‌های کاوش در داده همچون الگوریتم‌های تکاملی ژنتیک، نزدیک‌ترین همسایگی می‌توان الگوهای متفاوتی، متناسب با آن داده‌ها را استخراج نمود و به کشف الگوهای پنهان در پس آن اطلاعات رسید. بررسی این مسئله خود جای تحقیقات بیشتر بر روی داده‌های مناطق اقلیمی مختلف و مقایسه عوامل مؤثر در نتیجه خروجی آن‌ها دارد. در این تحقیق سه فاکتور مؤثر در پیدایش بیماری در نظر گرفته شده که می‌توان ورودی‌های در نظر گرفته شده را با در نظر گرفتن سایر عوامل ورودی همچون تأثیر متغیرهای دیگری چون

References

- Ovesen PG, Jensen DM, Damm P, Rasmussen S, Kesmodel US. Maternal and neonatal outcomes in pregnancies complicated by gestational diabetes. a nation-wide study. *J Matern Fetal Neonatal Med* 2015;28(14):1720-4.
- Coustan DR. Gestational diabetes mellitus. *Clinical Chemistry* 2013;59(9):1310-21.
- Kim C, Cheng YJ, Beckles GL. Cardiovascular disease risk profiles in women with histories of gestational diabetes but without current diabetes. *Obstet Gynecol* 2008;112(4):875-83.
- Lingaraj H, Devadass R, Gopi V, Palanisamy K. Prediction of diabetes mellitus using data mining techniques: a review. *Journal of Bioinformatics & Cheminformatics* 2015;1(1):1-3.
- Class T. Diabetes and pregnancy. *Diabetologia Croatica*. 2002;31:3.
- Chu SY, Callaghan WM, Kim SY, Schmid CH, Lau J, England LJ, et al. Maternal obesity and risk of gestational diabetes mellitus. *Diabetes Care* 2007;30(8):2070-6.
- O'Leary JA. *Shoulder Dystocia and Birth Injury: Prevention and Treatment*. 3th ed. USA: Humana Press; Softcover reprint of hardcover; 2009.
- Stafford GC, Kelley PE, Syka JEP, Reynolds WE, Todd JF. Recent improvements in and analytical applications of advanced ion trap technology. *International Journal of Mass Spectrometry and Ion Processes* 1984;60(1):85-98.
- Vijayarani S, Sudha S. Disease prediction in data mining technique—a survey. *International Journal of Computer Applications & Information Technology* 2013;2(1):17-21.
- Lakshmi KV, Padmavathamma M. Modeling an Expert System for Diagnosis of Gestational Diabetes Mellitus Based On Risk Factors. *Journal of Computer Engineering* 2013;8(3):29-32.
- Smitha V. An expert system for diabetes diagnosis. [dissertation]. India: Christ University Bangalore; 2010.
- Zarkogianni K, Vazeou A, Mougiakakou SG, Proutzou A, Nikita KS. An insulin infusion advisory system based on autotuning nonlinear model-predictive control. *IEEE Trans Biomed Eng* 2011;58(9):2467-77.
- Jaafar SFB, Ali DM. Diabetes mellitus forecast using artificial neural network (ANN). *Asian*

Conference on Sensors and the International Conference on New Techniques in Pharmaceutical and Biomedical Research; 2005 Sep 5-7; Kuala Lumpur, Malaysia, Malaysia: IEEE; 2005.

14. Sreedevi E, Vijaya Lakshmi K, Chaitanya Krishna E, Padmavathamma M. Modelling effective diagnosis of risk complications in gestational diabetes mellitus: an e-diabetic expert system for pregnant women. 4th International Conference on Digital Image Processing (ICDIP 2012); 2012 8 Jun 8; Kuala Lumpur, Malaysia: Proc. SPIE; 2012.

15. Gurumurthy S, Tripathy BK, Priya M. Study of Image Recognition Using Cellular Associated Artificial Neural Networks. Proceedings of the International MultiConference of Engineers and Computer Scientists; 2011. Mar 16-18; Hong Kong: IMECS; 2011.

16. Srivastava S, Tripathi K. Artificial neural network and non-linear regression: a comparative study. International Journal of Scientific and Research Publications 2011;2(12):1-5.

17. Sumathy M, Thirugnanam M, Kumar P, Jishnujit T, Kumar KR. Diagnosis of diabetes mellitus based on risk factors. International Journal of Computers and Applications 2010;10(4):1-4.

18. Nguyen HT, Ghevondian N, Nguyen ST, Jones TW. Detection of hypoglycemic episodes in children with type 1 diabetes using an optimal Bayesian neural network algorithm. Conf Proc IEEE Eng Med Biol Soc 2007;2007:3140-3.

19. Mirsharif M, Alborzi M. A fuzzy expert system & neuro-fuzzy system using soft computing for gestational diabetes mellitus diagnosis. International Journal of Information, Security and Systems Management 2014;3(1):249-52.

20. Su MC. Use of neural networks as medical diagnosis expert systems. Computers in Biology and Medicine 1994;24(6):419-29.

21. Jaafar SFB, Ali DM. Diabetes mellitus forecast using artificial neural network (ANN). Asian Conference on Sensors and the International Conference on New Techniques in Pharmaceutical and Biomedical Research; 2005 Sep 5-7; Kuala Lumpur, Malaysia, Malaysia IEEE; 2005.

22. Cunningham F, Leveno K, Bloom S, Hauth J, Rouse D, Spong C. Williams Obstetrics. 23th ed. New York: McGraw Hill ; 2010.

23. Hirose Y, Yamashita K, Hijjiya S. Back-propagation algorithm which varies the number of hidden units. Neural Networks 1991;4(1):61-6.

24. Sapra RL, Mehrotra S, Nundy S. Artificial neural networks: prediction of mortality/survival in gastroenterology. Current Medicine Research and Practice 2015;5(3):119-29.

25. Rumelhart D, Hinton G, Williams R. Learning internal representation by back propagation. Cambridge, MA, USA: Parallel distributed processing: explorations in the microstructure of cognition; 1986.

26. San PP, Ling SH, Nguyen HT. Intelligent detection of hypoglycemic episodes in children with type 1 diabetes using adaptive neural-fuzzy inference system. Conf Proc IEEE Eng Med Biol Soc 2012;2012:6325-8.

27. David SK, Saeb AT, Al Rubeaan K. Comparative analysis of data mining tools and classification techniques using weka in medical bioinformatics. Computer Engineering and Intelligent Systems 2013;4(13):28-38.

28. Hatonen K, Klemettinen M, Mannila H, Ronkainen P, Toivonen H, editors. Knowledge discovery from telecommunication network alarm databases. Proceedings of the Twelfth International Conference on Data Engineering; 1996 26 Feb-Mar 26-1; New Orleans, LA, USA, USA: IEEE; 1996.

29. Reichetzeder C, Dwi Putra SE, Pfab T, Slowinski T, Neuber C, Kleuser B, et al. Increased global placental DNA methylation levels are associated with gestational diabetes. Clin Epigenetics 2016;8:82.

30. Vijayan V, Ravikumar A. Study of data mining algorithms for prediction and diagnosis of diabetes mellitus. International Journal of Computer Applications 2014;95(17):12-6.

31. Devi MR, Shyla JM. Analysis of various data mining techniques to predict diabetes mellitus. International Journal of Applied Engineering Research. 2016;11(1):727-30.

Data Mining Approach based on Neural Network and Decision Tree Methods for the Early Diagnosis of Risk of Gestational Diabetes Mellitus

Mirsharif Maryam^{1*}, Rouhani Saeed²

• Received: 30 Apr, 2017

• Accepted: 20 Jan, 2017

Introduction: Nowadays, in this industrial modern world, the incidence of chronic diseases has been significantly increased. Gestational diabetes mellitus is one of the major health problems that if not treated, it will cause serious complications for mother and her child. The purpose of this research was to find ways for determining the risk of gestational diabetes mellitus and making early diagnosis to prevent it in the initial stages of pregnancy.

Methods: This applied-survey research used two approaches of neural network and decision tree in experimental analysis of data and prediction. The extracted data were normalized and analyzed through Matlab software.

Results: The results showed that data-based method is effective in improving the accuracy of prediction and has good performance in discovering implied knowledge and diagnosis of hidden relationships among data. In both methods, decision errors were acceptable and very close to each other.

Conclusion: Based on the obtained results, data mining methods can be used in health centers for less familiar diseases in order to achieve on-time diagnosis, patient management and to decrease treatment costs.

Keywords: Data mining, Artificial Neural Network, Decision Tree, Gestational diabetes mellitus, Diagnosis

• **Citation:** Mirsharif M, Rouhani S. Data Mining Approach based on Neural Network and Decision Tree Methods for the Early Diagnosis of Risk of Gestational Diabetes Mellitus. *Journal of Health and Biomedical Informatics* 2017; 4(1): 59-68.

1. MSc in Information Technology Management, Tehran University of Science and Research, Tehran, Iran.

2. Ph.D. in Systems Engineering, Assistant Professor, Faculty of Management, University of Tehran, Tehran, Iran

*Correspondence: University of Science and Research, Tehran, Iran.

• Tel: 09369566433

• Email: mm.mirsharif@gmail.com