



ارائه مدل پیش‌بینی تشخیص عوامل ناباروری؛ با استفاده از الگوریتم‌های داده کاوی

سمیرا درمحمدی^۱ / سمیه علیزاده^۲ / محسن اصغری^۳ / مریم شامی^۴

چکیده

مقدمه: حدود ۱۰-۱۵ درصد از زوجین نابارور هستند. ناباروری علل متفاوتی دارد و تشخیص روش درمان بیماران بر اساس نوع عامل ناباروری آن‌ها انجام می‌شود. در این تحقیق مدلی ارائه شده است که بر اساس ویژگی‌های اولیه و نتایج آزمایشات ساده علل ناباروری افراد را پیش‌بینی می‌کند که می‌تواند به پزشکان در تشخیص زودهنگام علت ناباروری و تصمیم‌گیری بهینه کمک کند.

روش کار: داده‌های این تحقیق برگرفته از داده‌های ناباروری بیمارستان صارم تهران می‌باشد. در این تحقیق از روش‌های داده‌کاوی استفاده شده است. ابتدا روش خوشه‌بندی k-means و سپس روش‌های دسته‌بندی ماشین بردار پشتیبان (SVM: Support Vector Machine) و شبکه‌های عصبی مصنوعی به منظور پیش‌بینی نوع علل ناباروری، اجرا و نتایج دو الگوریتم دسته‌بندی با هم مقایسه شدند. همچنین برای تحلیل داده‌ها و اجرای الگوریتم‌های بخش مدل، از نرم‌افزار SPSS Clementine 12.0 استفاده شده است.

یافته‌ها: در بخش خوشه‌بندی بر اساس الگوریتم K-means داده‌ها به پنج خوشه تقسیم شدند. در هر گروه یک یا چند علت ناباروری مشاهده شد. در ادامه و با اجرای الگوریتم‌های دسته‌بندی SVM و شبکه عصبی مصنوعی، مشخص شد که الگوریتم SVM با نوع کرنل چندجمله‌ای بالاترین کارایی را به دست آورد.

نتیجه‌گیری: انجام این تحقیق علاوه بر اینکه منجر به شناخت بهتر ویژگی‌های بیماران ناباروری شد، می‌تواند زمینه‌ای برای انجام تحقیقات آتی باشد. از آنجائی که با تشخیص علل ناباروری افراد قبل از مراحل ثانویه و آزمایشات سنگین، به مقدار قابل توجهی در هزینه و زمان صرفه‌جویی و از اثرات جسمی که بر بیماران می‌گذارد کاسته خواهد شد، می‌توان در مطالعات آینده با استفاده از نتایج این تحقیق سیستمی را جهت اجرای این مدل پیاده‌سازی نمود.

کلید واژه‌ها: ناباروری، مدل، داده کاوی، k-means، ماشین بردار پشتیبان، شبکه‌های عصبی مصنوعی

• وصول مقاله: ۹۲/۱۱/۲۱ • اصلاح نهایی: ۹۳/۲/۲۴ • پذیرش نهایی: ۹۳/۴/۴

۱. دانشجوی کارشناسی ارشد فناوری اطلاعات، دانشکده مهندسی صنایع، دانشگاه خواجه نصیر الدین طوسی، تهران، ایران، نویسنده مسئول (samira.dormohammadi@ymail.com)

۲. استادیار گروه مهندسی فناوری اطلاعات، دانشکده مهندسی صنایع، دانشگاه خواجه نصیر الدین طوسی، تهران، ایران

۳. کارشناسی ارشد، فناوری اطلاعات، دانشکده مهندسی صنایع، دانشگاه خواجه نصیر الدین طوسی، تهران، ایران

۴. کلینیک، ریاست، بیمارستان صارم، تهران، ایران

بارداری استفاده کنند» [۳]. حدود ۱۵-۱۰ درصد درصدا از زوجین نابارور هستند [۴]. دلایل گوناگونی برای ناباروری ذکر شده است از جمله علایمی مانند: فاکتور مردانه، فاکتور لوله‌ای، فاکتور دهانه رحم، علایم تخمک گذاری و علایم نامشخص. بر اساس تشخیص علت ناباروری که طی مراحل نسبتاً طولانی صورت می‌گیرد، یکی از روش‌های کمک باروری برای درمان انتخاب می‌شود. این روش‌ها شامل تلقیح داخل رحمی (IUI: Intra Uterine Insemination)، تلقیح اسپرم داخل سیتوپلاسم تخمک (ICSI: Intra Cytoplasmic Sperm Injection)، لقاح آزمایشگاهی (IVF: In Vitro Fertilization) و غیره است [۳].

در فرایند تشخیص علت ناباروری جلسه اول شامل گرفتن شرح حال در مورد تاریخچه قاعدگی، سابقه حاملگی قبلی، عمل جراحی، ابتلا به بیماری‌های زنان، مصرف داروها و درمان‌های قبلی نازایی می‌باشد. سایر بررسی‌ها شامل آزمایشات هورمونی به منظور بررسی عملکرد تخمدان‌ها است. این آزمایشات شامل تست‌های تیروئیدی، AMH، تستسترون، پرولاکتین، استروژن، LH و FSH می‌باشد. آزمایشات سرمی شامل HBSA، HBSAb، HBCAb، HCVAb، HIV، Rubella، CMV، Toxoplasma و قندخون و انسولین است. در اکثر زوج‌ها تست بعد از مقاربت (PCT: Post Coital Test) به منظور بررسی سلامت نطفه مرد و یا اثر ترشحات دهانه رحم زن بر نطفه‌های مرد انجام می‌شود. انجام یک سونوگرافی واژینال نیز در روزهای ۱۴-۱۲ جهت بررسی تخمک‌گذاری غالباً هم‌زمان با PCT توصیه می‌گردد. هیستروسالپینگوگرافی (HSG) نیز جهت بررسی باز بودن لوله‌ها و اطلاع از وضعیت حفره رحم توصیه می‌شود. نهایتاً پس از تکمیل مراحل فوق‌الذکر با توجه به نتیجه بررسی‌ها، درمان مورد نظر انتخاب و آماده‌سازی بیمار جهت پروتکل انتخابی آغاز می‌گردد.

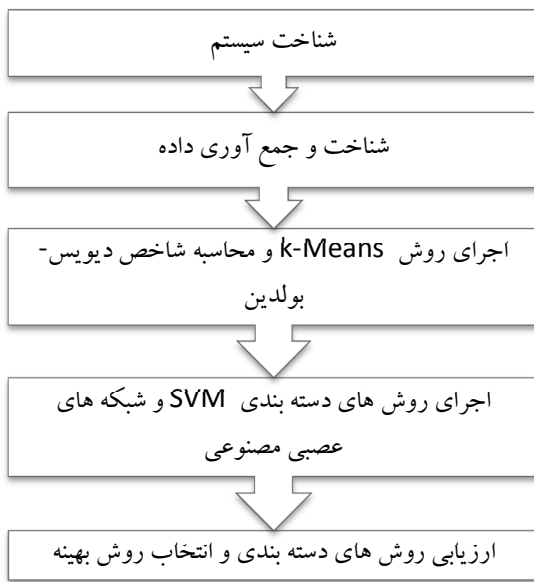
با توجه به مطالعات بررسی شده در زمینه داده‌های کاوی در ناباروری اشرفی کاخکی و همکاران [۵]، شیانگ گو و همکاران [۶]، اویار و همکاران [۸،۷]، والد [۹]، ونکات و

مقدمه

دستیابی به اینترنت و تکنیک‌های پیشرفته ذخیره اطلاعات، منابع و سازمان‌ها را با حجم زیادی از اطلاعات روبه‌رو کرده است، بدون اینکه دانش خاصی از آن‌ها دریافت شود. با رشد نمایی حجم اطلاعات تغییر ماهیت داده‌ها و زیاد شدن ابعاد داده (ویژگی‌های هر رکورد) ابزارهای تحلیل داده سنتی قادر به تحلیل داده‌ها نیستند. از طرفی سؤالاتی که در هر سازمان باید به آن‌ها پاسخ داده شود، نیاز به توسعه ابزارهای جدید را بیشتر کرده است. به عنوان راه حل این مسئله، تکنیک‌ها و ابزارهای داده‌کاوی توسعه پیدا کردند [۱].

«داده کاوی فرآیند کشف دانش مطلوب از مقدار بزرگی از داده است که در پایگاه داده، انبار داده و دیگر مخازن اطلاعات ذخیره شده است» [۲]. الگوریتم‌های داده‌کاوی به دو دسته کلی نظارتی و غیرنظارتی و یا پیش‌بینی و توصیفی تقسیم شده است. در الگوریتم‌های پیش‌بینی، هدف پیش‌بینی یک ویژگی خاص بر مبنای ویژگی‌های دیگر است. ویژگی پیش‌بینی شونده متغیر وابسته و بقیه متغیرها متغیر مستقل نامیده می‌شوند، اما در الگوریتم‌های توصیفی هدف استخراج الگو از داده‌ها است که نیاز به تحلیل نتایج دارد. الگوریتم خوشه بندی یکی از زیرمجموعه‌های الگوریتم‌های توصیفی است. در خوشه بندی، داده‌ها به گروه‌هایی به نام خوشه تقسیم می‌شوند به طوری که اعضاء یک خوشه بیشترین شباهت را به هم داشته باشند و اعضاء خوشه‌های مختلف کمترین شباهت را به هم داشته باشند [۱،۲].

از داده‌کاوی در زمینه‌های مختلفی از جمله مدیریت ارتباط با مشتری و تحلیل رفتار مشتریان، پیش‌بینی ارزش سهام، دسته بندی سبد محصول، پزشکی و هر جایی که حجم زیادی از داده موجود است استفاده می‌شود [۱]. یکی از شاخه‌های پزشکی که توجه محققان داده‌کاوی را به خود جلب کرده است، ناباروری می‌باشد. ناباروری عبارت است از «عدم وقوع بارداری تا یک سال پس از اینکه زوجین تصمیم به بچه‌دار شدن می‌گیرند، بدون اینکه از روش‌های پیش‌گیری از



شکل ۱: مراحل اجرای تحقیق حاضر

الف) شناخت سیستم

در این مرحله به شناخت زمینه مورد مطالعه پرداخته شده است. مواردی همچون شناخت داده‌ها، هدف از انجام این تحقیق، نیازمندی‌های مورد نیاز محیط تحقیق، مشکلات موجود در سیستم و موارد دیگر در این مرحله جای می‌گیرد.

ب) شناخت و جمع‌آوری داده‌ها

مجموعه داده این تحقیق شامل ۶۴۶ رکورد از اطلاعات بیماران بیمارستان صارم تهران می‌باشد که بین سال‌های ۱۳۸۵ تا ۱۳۹۰ جمع‌آوری شده است. در این مرحله عملیات پیش‌پردازش و پاکسازی داده بر روی آن‌ها جهت استفاده در ابزار تحلیل صورت گرفته است [۱۸، ۱۷]. مرحله پیش‌پردازش به منظور بهبود داده‌ها انجام می‌شود [۱] که شامل انجام فرایندهایی از قبیل: تصحیح و یا حذف داده‌های بدون مقدار، تعیین محدوده مجاز و تصحیح مقادیر غیرمجاز، انجام محاسبات مجدد برای برخی از ویژگی‌ها و تبدیل آن‌ها به ویژگی‌های دیگر است [۱۸، ۱۷]. هر رکورد از پایگاه داده، ۱۵ ویژگی دارد. مشخصات مربوط به هر متغیر و شرح مختصری از آن در جدول (۱) نشان داده شده است.

همکاران [۱۰]، جورسیکا و همکاران [۱۱] و کافمن و همکاران [۱۲] به پیش‌بینی نتیجه عمل کمک باروری IVF پرداختند. در بین این مطالعات مدل پیشنهادی در اشرافی کاخکی و همکاران [۵] بیشترین صحت (۹۳ درصد) را به دست آورده است. میلسکی و همکاران به دسته بندی داده های بیماران IVF و ICSI پرداختند [۱۳]. مورالز و همکاران نیز با استفاده از روش دسته بندی بیزین به پیش‌بینی بهترین جنین آزمایشگاهی پرداختند [۱۴]. در مطالعات دیگر، میکاس و همکاران به بررسی یکی از علل ناباروری (مشکل آروسپریمیا که یکی از انواع مشکل در علت فاکتور مردانه در ناباروری است) [۱۵] و زروسکی و همکاران به مشکل حذف کروموزومی Y پرداخته است [۱۶] و دیگر عوامل ناباروری در هیچ کدام از مطالعاتی که در زمینه داده کاوی در ناباروری مورد مطالعه قرار گرفته، بررسی نشده است. بنابر وجود چنین خلأیی، در این تحقیق قرار شد تا با استفاده از ویژگی‌های اولیه و نتایج آزمایشات ساده بیماران ناباروری، مدلی ارائه شود که ابتدا این بیماران را به گروه‌های مختلف تقسیم نموده و سپس بر اساس یک روش پیش‌بینی، دسته بیمار جدید را تشخیص داده و علل ناباروری احتمالی را برای آن پیش‌بینی نماید. در صورتی که بتوان پیش‌بینی نمود که علت ناباروری افراد چیست، می‌توان با حذف یک مرحله و یا ترتیب درست انجام آزمایشات، در هزینه و زمان برای بیمار صرفه‌جویی و فرایند درمان بیماران را بهینه نمود. همچنین نتایج این تحقیق می‌تواند در فرایند تصمیم‌گیری در مورد نوع روش درمانی که بر اساس نوع علت ناباروری صورت می‌گیرد به پزشک کمک نماید.

روش کار

مطالعه حاضر داده‌محور می‌باشد که در آن تحقیقی با استفاده از روش‌های داده‌کاوی بر روی داده‌های ناباروری صورت گرفته است. مراحل مختلف اجرای تحقیق حاضر در شکل (۱) نشان داده شده است.

جدول ۱: متغیرهای تحقیق

نام ویژگی	شرح	نوع	بازه مقادیر
BMI Status	شاخص توده بدنی (تناسب قد و وزن زن)	اسمی	لاغر، نرمال، اضافه وزن و چاق
OF	وجود مشکل تخمک گذاری	اسمی	مثبت، منفی
TF	وجود مشکل لوله‌ای	اسمی	مثبت، منفی
MF	وجود مشکل مردانه	اسمی	مثبت، منفی
Infertility	نوع ناباروری	اسمی	اولیه، ثانویه
Kind Of Protocol	نوع پروتکل	اسمی	Poor, Long, Short
Embryo No Transform	تعداد جنین منتقل شده	عددی	بازه [۰-۷]
No Of Ovum	تعداد تخمک‌ها	عددی	بازه [۰-۴۶] میلی‌متر
Duration	طول مدت ناباروری	عددی	بازه [۱-۳۰] میلی‌متر
Female Hormonal Test FSH	میزان هورمون FSH زن	عددی	بازه [۰/۱-۷۶]
Female Hormonal Test LH	میزان هورمون LH زن	عددی	بازه [۰/۱-۱۰۱]
Female Hormonal Test Estradiol	میزان هورمون استروژن زن	عددی	بازه [۰/۶-۶۷۳]
Age	سن زن	عددی	بازه [۲۰-۵۰]
No Follicle	تعداد فولیکول	عددی	بازه [۱-۳۸] میلی‌متر
Thickness	ضخامت آندومتر (پوشش داخلی رحم)	عددی	بازه [۴-۱۷] میلی‌متر

(Bouldin)، سیلوهونت (Silhouette) و غیره تعریف شده است. انتخاب k بهینه در شاخص دیویس-بولدین بر مبنای اصل «کمترین شباهت بین خوشه‌ای و بیشترین شباهت درون خوشه‌ای» استوار است. در این تحقیق شاخص دیویس-بولدین به دلیل گستردگی استفاده به کار گرفته شده است.

در معادلات زیر V مرکز خوشه، DB مقدار نهایی شاخص دیویس بولدین، d نشان دهنده فاصله خوشه‌ها از یکدیگر، S نشان دهنده پراکندگی داخل خوشه، q و t یک عدد صحیح، A_i مجموعه‌ای از رکوردهایی است که در هر مرحله در خوشه قرار می‌گیرد، $|A_i|$ تعداد عناصر مجموعه A_i و C نشان دهنده تعداد خوشه‌ها در هر مرحله از محاسبه شاخص است [۱۹]. بنابراین داریم:

(معادله ۱) [۱۹]

$$S_{i,q} = \left(\frac{1}{|A_i|} \sum_{x \in A_i} \|x - V_i\|_2^q \right)^{1/q}$$

(معادله ۲) [۱۹]

$$d_{ij,s} = \left\{ \sum_{s=1}^p |v_{si} - v_{sj}|^t \right\}^{1/t} = \left\| \underline{v}_i - \underline{v}_j \right\|_t$$

همانطور که در جدول (۱) ملاحظه نمودید برخی از ویژگی‌ها که اغلب آن‌ها اسمی هستند مقادیر مشخصی دارند.

ج) اجرای روش k -means و محاسبه شاخص دیویس-بولدین

از معروف‌ترین الگوریتم‌های خوشه بندی k -means است. در این الگوریتم ابتدا k مقدار از اشیاء به عنوان مراکز اولیه خوشه‌ها انتخاب می‌شود و سپس فاصله هر شیء با این مراکز خوشه محاسبه می‌شود و هر شیء به خوشه‌ای اختصاص می‌یابد که فاصله آن با مرکز خوشه کمترین باشد. سپس مقدار مرکز جدید خوشه محاسبه می‌شود. این فرآیند تا زمانی که تابع معیار همگرا شود ادامه می‌یابد [۲].

در این تحقیق از نرم‌افزار داده‌کاوی SPSS Clementine 12.0 استفاده شده است. یکی از مهم‌ترین معایب الگوریتم k -means، انتخاب k است که از ابتدا مشخص نیست. ابتدا باید این الگوریتم را برای k های مختلف محاسبه کنیم و سپس با استفاده از یک شاخص بررسی کنیم که کدام k خوشه بندی بهینه را به ما می‌دهد [۲]. در منابع مختلف شاخص‌های متفاوتی از جمله شاخص دان (Dunn)، دیویس-بولدین (Davies-

[معادله ۳] [۱۹]

$$R_{i,qt} = \max_{j \in c, j=1} \left\{ \frac{S_{i,q} + S_{j,q}}{d_{ij,s}} \right\}$$

و معادله نهایی شاخص دیویس-بولدین برای تعداد c خوشه به صورت زیر است:

$$DB(c) = \frac{1}{c} \sum_{i=1}^c R_{i,q} \quad [معادله ۴] [۱۹]$$

تعداد خوشه بهینه برابر کمترین مقدار $DB(c)$ است. الگوریتم k -means با k مقادیر k ، از $k=3$ تا $k=8$ بر روی مجموعه داده حاصل اجرا شد. بقیه مقادیر الگوریتم، همان مقادیر پیش فرض است و تمام ویژگی‌های نمونه‌های مجموعه داده در خوشه‌بندی به عنوان ورودی انتخاب شدند.

د) اجرای روش‌های دسته‌بندی SVM و شبکه عصبی مصنوعی الگوریتم SVM یکی از الگوریتم‌هایی است که در زیرگروه الگوریتم‌های پیش‌بینی قرار می‌گیرد و مبنای آن از تئوری یادگیری آماری است. SVM با استفاده از یک خط به نام مرز تصمیم، نمونه‌های کلاس‌های مختلف را از هم جدا می‌کند. این مرز تصمیم بردار پشتیبان نامیده می‌شود. هر مرز تصمیم به دو ابرصفحه (hyper plan) محدود می‌شود که فاصله آن‌ها از مرز تصمیم یکسان است و فاصله بین این دو ابرصفحه حاشیه دسته‌کننده است. هدف دسته‌کننده SVM پیدا کردن یک مرز تصمیم با حاشیه دسته‌کننده حداکثر است. مرز تصمیم می‌تواند خطی و یا غیرخطی باشد. این الگوریتم برای داده‌هایی با ابعاد بالا خوب عمل می‌کند [۱].

یکی دیگر از الگوریتم‌های پیش‌بینی شبکه‌های عصبی مصنوعی است. این الگوریتم قصد دارد برخی از عملکردهای ساده مغز انسان را شبیه‌سازی کند. هر گره در شبکه عصبی مانند یک نرون در مغز انسان است که شبکه اتصال این نرون‌ها وظایف یادگیری پیچیده‌ای را انجام می‌دهند. فرایند تحلیل داده‌ها در شبکه عصبی مانند یک جعبه سیاه است. شبکه‌های عصبی در تخمین و پیش‌بینی مثلاً در تخمین قیمت سهام در ماه بعد، بسیار کاربرد دارد [۲].

پس از تعیین گروه‌های مختلف، روش دسته‌بندی SVM و شبکه‌های عصبی برای پیش‌بینی گروه مربوط به نمونه جدید

اجرا گردید و نتیجه آن‌ها با هم مقایسه شد. در روش‌های دسته‌بندی، یک فیلد باید به عنوان فیلد خروجی انتخاب شود که در اینجا، خوشه‌ها به عنوان فیلد خروجی انتخاب شد و فیلدهای وضعیت BMI، نوع ناباروری، ضخامت آندومتر، سن، تعداد فولیکول، تست هورمونی LH زن، تست هورمونی FSH زن، تست هورمونی استروژن زن، طول مدت ناباروری و تعداد تخمک به عنوان فیلدهای ورودی در نظر گرفته شدند.

در الگوریتم SVM، در نرم‌افزار SPSS Clementine، کرنل، نوع خط مرز تصمیم را نشان می‌دهد و چهار مقدار مختلف می‌پذیرد و با تغییر مقدار این ویژگی نتایج متفاوت حاصل می‌گردد. این مقادیر شامل خطی (Linear)، چندجمله‌ای (Polynomial)، حلقوی (Sigmoid) و RBF است که همه این چهار روش در بخش دسته‌بندی این تحقیق اجرا شد و با هم مقایسه گردید. پارامترهای دیگر الگوریتم SVM شامل Regression Precision، Regularization Parameter (epsilon)، RBF gamma، Bias و Degree است که در این تحقیق به ترتیب دارای مقادیر ۰/۰۵، ۰/۰۱، ۰/۰۵، ۰/۰۱ و ۴ هستند. الگوریتم شبکه عصبی مصنوعی نیز با مقادیر پیش فرض اجرا گردید.

ه) ارزیابی روش‌های دسته‌بندی و انتخاب روش بهینه در این مرحله ارزیابی مدل صورت می‌گیرد. صحت، یکی از معیارهای ارزیابی مدل‌های دسته‌بندی است که مقدار آن برابر درصد مشاهدات مجموعه آموزشی است که توسط روش مورد استفاده، به درستی دسته‌بندی شده است. ماتریس اغتشاش (Confusion) یکی از ابزارهای مفید برای ارزیابی عملکرد روش‌های دسته‌بندی است. اگر تعداد دسته‌های موجود m باشد ماتریس، اغتشاش جدولی با اندازه $m \times m$ است. اگر i شماره سطر باشد و j شماره ستون باشد عنصر C_{ij} تعداد مشاهداتی از دسته i است که توسط الگوریتم دسته‌بندی تشخیص داده شده است [۲۰]. معیارهای دیگر برای ارزیابی عملکرد الگوریتم دسته‌بندی، حساسیت (Sensitivity)، شفافیت (Specificity)، دقت (Precision) و صحت است [۲۱] که این مقادیر در ذیل تعریف شده اند:

$$\text{حساسیت} = \frac{\text{تعداد داده های کلاس مثبت که درست دسته بندی شدند}}{\text{تعداد کل داده های کلاس مثبت}} \quad (\text{معادله ۵})$$

$$\text{شفافیت} = \frac{\text{تعداد داده های کلاس منفی که درست دسته بندی شدند}}{\text{تعداد کل داده های کلاس منفی}} \quad (\text{معادله ۶})$$

$$\text{دقت} = \frac{\text{تعداد داده های کلاس مثبت که درست دسته بندی شدند} + \text{تعداد داده های کلاس منفی که درست دسته بندی شدند}}{\text{تعداد داده های کلاس مثبت که درست دسته بندی شدند} + \text{تعداد داده های کلاس منفی که درست دسته بندی شدند}} \quad (\text{معادله ۷})$$

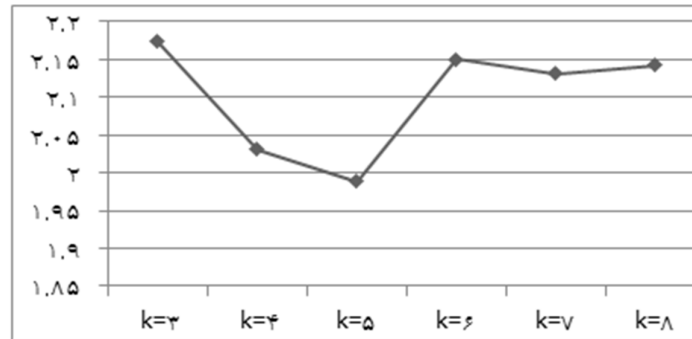
$$\text{صحت} = \frac{\text{تعداد داده های کلاس مثبت}}{\text{تعداد کل داده ها}} * \text{حساسیت} + \frac{\text{تعداد داده های کلاس منفی}}{\text{تعداد کل داده ها}} * \text{شفافیت} \quad (\text{معادله ۸})$$

داده‌های اعتبارسنجی (۱۰ درصد) برای بررسی ارزش و اعتبار تقسیم شدند.

یافته‌ها

پس از اجرای الگوریتم k-means، تعداد خوشه بهینه محاسبه شد که برابر مینیمم شاخص دیویس بولدین (شکل ۲) یعنی $k=5$ است.

در فرمول‌های ۵، ۶، ۷ و ۸ منظور از کلاس مثبت کلاسی است که شاخص برای آن محاسبه می‌شود (یکی از پنج خوشه) و کلاس منفی مابقی کلاس‌ها به جز کلاس انتخاب شده است [۲]. قبل از اجرای دسته‌بندی، داده‌ها به سه دسته داده‌های آموزش (۶۵ درصد) برای ساخت روش دسته بندی و داده‌های تست (۲۵ درصد) برای آزمایش دسته بندی و



شکل ۲: مقدر شاخص دیویس بولدین برای مجموعه داده با نتایج منفی

مردانه دارند و ۱/۳ درصد از آنان هیچ علت شناخته شده‌ای ندارند و جز بیماران با علل ناشناخته هستند.

خوشه دو: بیماران این خوشه بالاترین میانگین سنی (۳۳/۶) را دارند و افرادی که مشکل عامل مردانه نداشتند (حدود ۲۱ درصد) در این خوشه نسبت به همه خوشه‌ها بیشتر است. بیشترین میزان (۲۲/۶ درصد) استفاده از پروتکل Long در این خوشه است. حدود ۶۰ درصد وزن غیرنرمال دارند. همه افراد آن مشکل لوله‌ای (TF) دارند، همچنین حدود ۲۳ درصد علت تخمک گذاری دارند، نابراین وجود علت تخمک گذاری نسبت به علل دیگر با ابهام بیشتری روبرو است.

نتایج حاصل از خوشه‌بندی k-means (خوشه‌ها) در جدول (۲) نشان داده شده است. برای متغیرهایی که مقادیر آن‌ها پیوسته و در بازه‌ی مشخصی قرار دارد، میانگین مقادیر و انحراف معیار که نشان دهنده مقدار پراکندگی مقادیر یک متغیر حول مقدار میانگین است، محاسبه شده است.

خوشه یک: تعداد بیماران (۱۱۰ نفر) در این خوشه نسبت به سایر خوشه‌ها کمتر است و بیماران این خوشه پایین‌ترین میانگین سنی (۳۰/۳) و کمترین طول مدت ناباروری (۶/۲) را دارند. ۹۸/۷ درصد از بیماران این خوشه فقط علت ناباروری

جدول ۲: خوشه‌های حاصل از اجرای الگوریتم k-means بر روی داده‌های بیماران با نتایج بارداری منفی

خوشه ۵	خوشه ۴	خوشه ۳	خوشه ۲	خوشه ۱	نام ویژگی
۱۳۴	۱۲۰	۱۵۸	۱۲۴	۱۱۰	تعداد رکورد
۳۱/۸	۳۲/۲	۳۱/۶	۳۳/۶	۳۰/۳	میانگین سن
۶/۱	۵/۶	۶/۶	۵/۱	۵	انحراف معیار سن
۷/۵	۸/۸	۸	۷/۸	۶/۲	میانگین مدت ناباروری
۵/۳	۵/۴	۴/۹	۵/۵	۵	انحراف معیار مدت ناباروری
۵/۸	۴/۷	۵/۹	۵	۴/۵	میانگین هورمون FSH
۴/۲	۴/۷	۶/۶	۳/۵	۲/۶	انحراف معیار هورمون FSH
۹/۹	۷/۸	۱۰/۴	۷/۲	۷/۵	میانگین هورمون LH
۱۲/۶	۸/۱	۱۱/۴	۷/۴	۴/۸	انحراف معیار هورمون LH
۸۶	۱۰۳/۱	۷۷/۹	۸۴/۹	۸۴/۱	میانگین هورمون Estradiol
۶۷/۱	۹۲/۸	۷۷/۶	۸۱/۶	۶۱/۷	انحراف معیار هورمون Estradiol
۸۷	۱۰/۶۷	۱۲/۷	۹/۱	۱۲/۹	میانگین تعداد تخمک
۶/۵	۷	۸/۴۶	۶/۳	۶/۵	انحراف معیار تعداد تخمک
۸/۱	۹/۲	۱۱/۳	۷/۷	۱۰/۶	میانگین تعداد فولیکول
۴/۶	۵/۵	۶	۴/۶	۵	انحراف معیار تعداد فولیکول
۹/۳۲	۹/۶۵	۹/۰۹	۹/۲	۹/۱	میانگین ضخامت آندومتر
۱/۹	۱/۸	۱/۵	۱/۸	۱/۶	انحراف معیار ضخامت آندومتر
%۱/۱۵	۰	%۰/۱۶	%۱/۱۶	%۰/۹	BMI > ۱۸/۹ لاغر
%۸۶/۶	۰	%۵۳/۸	%۳۸/۷	%۸۸/۲	BMI < ۲۴/۹ < ۱۹ نرمال
۰	%۱۰۰	%۲۹/۱	%۴۶	۰	BMI < ۲۹/۹ < ۲۵ اضافه وزن
%۱۱/۹	۰	%۱۶/۴	%۱۳/۷	%۱۰/۹	BMI < ۳۰ چاق
%۸۵/۸	%۸۷/۵	%۷۵/۳	%۶۲/۹	%۸۲/۷	اولیه نوع ناباروری
%۱۴/۲	%۱۲/۵	%۲۴/۷	%۳۷/۱	%۱۷/۳	ثانویه نوع ناباروری
%۷۷/۶	%۸۷/۵	%۸۵/۴	%۷۱	%۹۰/۹	Short نوع پروتکل
%۱۴/۲	%۶/۷	%۱۱/۴	%۲۲/۶	%۶/۴	Long نوع پروتکل
%۸/۲	%۵/۸	%۳/۲	%۶/۵	%۲/۷	Poor نوع پروتکل
%۰/۸	۰	۰	۰	۰	۰
%۱۴/۹	%۶/۷	%۳/۸	%۱۴/۵	۰	۱
%۳۰/۶	%۱۸/۳	%۱۵/۸	%۱۷/۷	۰	۲
%۳۵/۰۷	%۲۰	%۲۰/۳	%۲۱/۸	۰	۳
۰	%۳۹/۲	%۴۶/۸	%۳۸/۷	%۱۰۰	۴
%۱۶/۴	%۱۴/۲	%۱۱/۴	%۷/۳	۰	۵
%۱/۵	%۱/۷	%۱/۹	۰	۰	۶
%۰/۸	۰	۰	۰	۰	۷
۰	۰	%۱۰۰	%۲۳/۴	۰	مثبت OF
%۹۹/۲	%۹۹/۲	%۹۰/۵	%۷۹	%۹۸/۲	مثبت MF
۰	۰	۰	%۱۰۰	۰	مثبت TF

مشکل لوله‌های فالوپ و تخمک گذاری ندارند. این بیماران بالاترین میانگین هورمون استروژن را دارند. خوشه پنج: حدود ۸۶ درصد از افراد این خوشه وزن نرمال دارند و حدود ۹۹ درصد از آن‌ها نیز همانند خوشه (چهار) فقط مشکل مردانه دارند و حدود یک درصد جز افراد با علل ناشناخته هستند و هیچ کدام از این افراد مشکل لوله‌ای و تخمک گذاری ندارند. بعد از اجرای الگوریتم‌های دسته بندی صحت الگوریتم SVM برای انواع کرنل و شبکه عصبی در جدول (۳) نشان داده شده است.

خوشه سه: تعداد افراد این خوشه (۱۵۸ نفر) نسبت به سایر خوشه‌ها بیشتر است. حدود ۴۷ درصد از بیماران این خوشه وزن غیرنرمال دارند. این تنها خوشه‌ای است که همه بیماران آن مشکل تخمک گذاری دارند، اما هیچ یک از آن‌ها مشکل لوله‌های فالوپ ندارند. خوشه چهار: همه افراد این خوشه اضافه وزن دارند. حدود ۹۹ درصد از آن‌ها فقط مشکل عامل مردانه دارند و بقیه جز بیماران با علل ناباروری ناشناخته هستند و

جدول ۳: مقدار صحت الگوریتم‌های متفاوت

SVM (کرنل چندجمله‌ای) (درصد)	شبکه عصبی (درصد)	
۷۵/۷	۴۱/۸	داده‌های آموزشی
۷۶/۷	۴۸	داده‌های تست
۷۴	۴۲/۵	داده‌های اعتبارسنجی

صحت برای برچسب کلاس‌های مختلف برای روش بهینه در جدول (۴) نشان داده شده است.

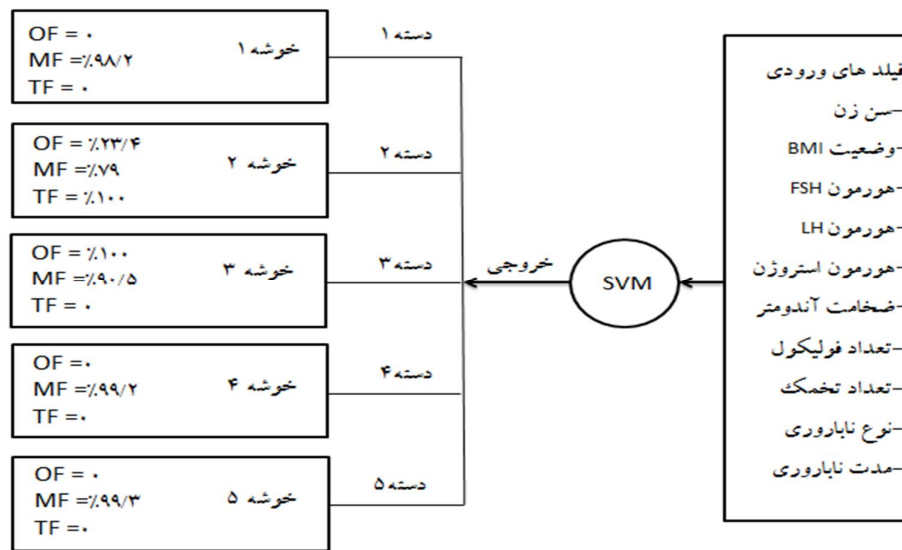
همانطور که ملاحظه می‌شود، الگوریتم SVM با کرنل چندجمله‌ای به عنوان الگوریتم دسته بندی با بیشترین صحت انتخاب شد. مقادیر معیارهای مختلف میزان

جدول ۴: مقادیر شاخص‌های مختلف برای الگوریتم SVM با تابع کرنل چندجمله‌ای

نام خوشه	حساسیت (درصد)	شفافیت (درصد)	دقت (درصد)	صحت (درصد)
خوشه ۱	۷۴/۵	۷۵/۹	۶۷/۲	۷۵/۷
خوشه ۲	۶۲/۹	۷۸/۷	۸۱/۳	۷۶/۴
خوشه ۳	۷۵/۳	۷۵/۸	۷۴/۸	۷۴/۶
خوشه ۴	۸۹/۲	۷۲/۶	۸۱	۷۶
خوشه ۵	۷۶/۹	۷۵/۴	۷۵/۲	۵۷/۷

خوشه‌های بدست آمده در بخش الگوریتم K-Means هستند قرار می‌دهد و در هر یک از آن‌ها احتمال وجود عوامل ناباروری مشخص شده است.

برای درک بهتر الگوریتم SVM، عملکرد آن به صورت شماتیک در شکل (۳) نشان داده شده است. با توجه به شکل، الگوریتم SVM بر اساس ویژگی‌های اولیه، بیمار را در یکی از دسته‌ها که همان



شکل ۳: عملکرد الگوریتم SVM

مشکل لوله‌ای دارد، با احتمال ۷۹ درصد مشکل مردانه دارد و با احتمال کم (۲۳ درصد) علت ناباروری بیمار عامل تخمک‌گذاری است. بنابراین در این حالت ابتدا می‌توان فاکتور لوله‌ای که احتمال وجود آن ۱۰۰ درصد است را بررسی کرد. نهایتاً اگر الگوریتم دسته‌بندی خوشه (سه) را برای نمونه بیمار جدید پیش‌بینی نمود، آنگاه می‌توانیم بگوئیم که این بیمار با احتمال ۱۰۰ درصد مشکل تخمک‌گذاری دارد و با احتمال ۹۷/۲ درصد مشکل مردانه دارد و مشکل لوله‌ای ندارد، بنابراین بررسی عامل تخمک‌گذاری در ابتدا منطقی‌تر به نظر می‌رسد.

با تشخیص زودهنگام علت ناباروری با استفاده از نتایج آزمایشات ساده اولیه می‌توان به روند درمان جهت داده و آن را بهینه نمود. چرا که روش‌هایی که برای تشخیص علل ناباروری بکار گرفته می‌شود، عوارض زیادی به همراه دارد، بنابراین انجام بیهوده آن برای بیمار عذاب‌آور است. به عنوان مثال، هیستروسالپنگوگرافی که می‌تواند برای زن عوارضی از قبیل درد، خونریزی و عوارض ناشی از تاباندن اشعه را داشته باشد [۲۱]. تصور اینکه حذف این مرحله چه اثرات جسمی و مادی مثبتی برای بیمار به همراه خواهد داشت، به دور از ذهن نیست.

تأثیر ویژگی‌ها در دسته‌بندی متفاوت است. SVM با کرنل چندجمله‌ای ویژگی‌های وضعیت BMI، تعداد تخمک، سن، تست استرادیول، تعداد فولیکول، نوع ناباروری، تست FSH را به ترتیب به عنوان مهم‌ترین ویژگی‌ها در دسته‌بندی معرفی کرد.

بحث و نتیجه‌گیری

در این تحقیق با استفاده از ترکیب الگوریتم‌های داده‌کاوی، علت ناباروری افراد با استفاده از ویژگی‌های اولیه آن‌ها پیش‌بینی شد. در بخش خوشه‌بندی، داده‌ها به پنج خوشه تقسیم شدند. علت ناباروری (ویژگی هدف) در هر یک از خوشه‌ها متفاوت بود. لازم به ذکر است که خوشه‌های (یک، چهار و پنج) نتایج مشابهی در علل ناباروری داشتند. در بخش پیش‌بینی، الگوریتم SVM با کرنل چندجمله‌ای بالاترین صحت را به دست آورد. برای داده‌های آموزشی میزان صحت ۷۵/۷ درصد، برای داده‌های تست ۷۶/۷ درصد و برای داده‌های اعتبارسنجی ۷۴ درصد حاصل شد. در صورتی که الگوریتم SVM بیماری با ویژگی خاص را در خوشه‌های (یک، چهار و پنج) قرار دهد، این بیمار فقط مشکل مردانه دارد. اگر در خوشه (دو) قرار دهد با احتمال ۱۰۰ درصد

References

1. Ghazanfari M, Alizadeh S, Teymourpour B. editors. [Data mining and knowledge discovery]. 2nd ed. Tehran: publication of university of science and technology ;2011.[Persian]
2. Han J, Kamber M . Data Mining: Concepts and Techniques. 2nd ed.San Francisco: Morgan Kufman Publisher; 2006.
3. Saremi A. [Introduction to Infertility]. 2nd ed. Tehran: Author; 2004. [Persian]
4. Esmailzadeh S, Farsi M, Bijani A. [Impact of morphology of sperm on the rate of fertility in Intra Uterine Insemination method]. Journal of Reproduction and Infertility 2007; 8(32): 205-212.[Persian]
5. Ashrafi Kakhki S, Malekara B, Rahati Quchani S, Khadem N . A model a Bayesian network for prediction of IVF success rate. 7th international symposium on advances in science and technology. Bandar-Abbas Iran. 2013 March.
6. Shiang Guh R, Chieh Jackson wu T, Ping Weng S. Integrating genetic algorithm and decision tree learning for assistance in predicting in vitro fertilization outcomes . Expert Systems with Applications 2011; 38: 4437-4449.
7. Uyar A, Bener A, Ciray H N, Bahceci M. ROC Based Evaluation and Comparison of Classifiers for IVF Implantation Prediction. Electronic Healthcare 2010; 27: 108-111.
8. Uyar A, Bener A, Ciray H N, Bahceci M. Handling the imbalance problem of IVF

طبق بررسی‌های انجام شده در زمینه داده‌کاوی در ناباروری، تحقیق میکاس و همکاران و همچنین زروسکی و همکاران در زمینه عوامل ناباروری است، که به گروه بندی یک نوع از بیماران با عوامل مردانه می‌پردازد. ولی به نظر می‌رسد تاکنون در هیچ مطالعه‌ای مدلی که علل ناباروری افراد را با استفاده از ابزار داده‌کاوی، پیش‌بینی کند ارائه نشده است و تحقیق حاضر اولین مطالعه در این زمینه است. نتایج حاصل از پیش‌بینی انجام شده در این تحقیق فقط می‌تواند به عنوان کمکی در تصمیم‌گیری استفاده شود و به هیچ عنوان نمی‌تواند جایگزین پزشک شود. همچنین نتایج این تحقیق وابسته به داده‌های بیمارستان صارم تهران می‌باشد. برای بررسی بیشتر در این زمینه می‌توان در مطالعات بعدی از داده‌های مراکز درمانی دیگر و الگوریتم‌های دیگر نیز استفاده کرده و نتایج را با هم مقایسه نمود.

تشکر و قدردانی

از زحمات جناب آقای دکتر صارمی و پرسنل محترم بیمارستان که در مراحل انجام این تحقیق همکاری لازم را داشتند، کمال تشکر و قدردانی را داریم.

9. implantation prediction. IAENG International Journal of Computer Science 2010 May; 37(2): 164-170.
10. Wald M. Computational models for prediction of IVF/ICSI outcomes with surgically retrieved spermatozoa. Reproductive Bio Medicine Online 2005 July; 11(3): 325-331.
11. Venkat G, Al-Nasser R, Jercovic S, Craft I. Prediction of success in IVF treatments using neural networks. Fertility and Sterility 2004 Sep; 82: s215.
12. Jurisica I, Mylopoulos J, Glasgow J, Shapiro H, Casper R. Case-based reasoning in IVF: prediction and knowledge mining. Artificial Intelligence in Medicine 1998; 12: 1-24.
13. Kaufman S J, Eastaugh JL, Snowden S, Smye SW, Sharma V. The application of neural networks in predicting the outcome of in-vitro fertilization. Human Reproduction 1997; 12(7): 1454-1457
14. Milewski R, Malinowski P, Milewski AJ, Ziniewicz P, Wolczynski S. Classification issue in IVF ICSI/ET data analysis. Studies in logic, grammar and rhetoric 2012; 29(42): 75-85.
15. Morales D A, Bengoetxea E, Larranaga P, Garcia M, Fresnada Y, Merino M. Bayesian classification for the selection of in vitro human embryos using morphological and clinical data. computer methods and programs in biomedicine 2008; 90: 104-116.
16. Mikos T, Pantazis K, Goulis D G, Maglaveras N, Bontis J N, papadimas J. The use of Data Mining in the categorization of patients with Azoospermia. HORMONES 2005 October; 4(4): 214-218.
17. Dzeroski S, Hristovski D, Peterlin B. Using Data Mining and OLAP to Discover Patterns in a Database of Patients with Y-Chromosome Deletions. Journal of the American Medical Informatics Association (AMIA) 2000; 7(Suppl): 215-219.
18. Hoseini M. [Developing a predictive model based on the Sarem hospital infertility data] [MSc thesis]. Tehran: K.N. Toosi University of technology; 2012. [Persian]
19. Aghabeigi N. [Providing a hybrid model for clustering on infertility data] [MSc thesis]. Tehran: K.N. Toosi University of technology; 2012. [Persian]
20. Abolmasum F, Alizadeh S, Asghari M. [Utilizing Data Mining Techniques for Investigating Factors Influencing the Failure of Intrauterine Insemination Infertility Treatment]. Journal of Tehran university of medical science 2013; 16(54) : 46-55. [Persian]
21. Ameri H, Alizade S, Barzegari A. [Knowledge Extraction of Diabetics' Data by Decision Tree Method]. Journal of Tehran university of medical science 2013 ; 16(53) : 58-72. [Persian]
22. Saghafi N, Farzaneh S. [Comparison of two methods of hystrosalpyngography and hydrososnohy sterosalpingography in evaluation of the uterine cavity and fallopian tubes in infertile women]. Journal of Iran's women, midwifery and infertility 2003; 6: 35-39. [Persian]



Proposing a prediction model for diagnosing Causes of Infertility by Data Mining Algorithms

Dormohammadi S¹/ Alizadeh S²/ Asghari M³/ Shami M⁴

Abstract

Introduction: About 10-15 percent of Iranian couples are infertile which is due to different causes determining particular diagnostic and treatment methods. In this study, the model presented is based on basic features and simple tests, helping physicians predict the causes of infertility

Methods: The data were taken from Sarem hospital infertility data bank by using data mining methods. First, K-means clustering was run; then, support vector machine and artificial neural network classification methods were used to predict the type of infertility, and finally, the results of two classification algorithms were compared. In addition, SPSS Clementine 12.0 was used to analyze the data and implement the algorithm in modeling part.

Results: In k-means clustering, the data were divided into five clusters. In each cluster, one or more causes of infertility were observed. Then, by applying SVM and artificial neural network classification algorithms, the SVM algorithm with a polynomial kernel appeared to have the maximum accuracy.

Conclusion: The findings of this study, could contribute to the understanding of the factors responsible for infertility and pave the way for future investigations. These findings can be used in future studies to develop a system for applying this model since by diagnosing the causes of infertility prior to secondary stages and before performing heavy tests, a considerable amount of time and cost will be saved, and physical burden on patient will be decreased,

Keywords: Infertility, Model, Data Mining, k-means, Support Vector Machine, Artificial Neural Network.

• Received: 10/Feb/2014 • Modified: 14/May/2014 • Accepted: 25/June/2014

1. MSc Student of Information Technology, Faculty of Industrial Engineering, K.N Toosi University of Technology, Tehran, Iran; Corresponding Author (samira.dormohammadi@ymail.com)
2. Assistant Professor of Information Technology Department, Faculty of Industrial Engineering, KN Toosi University Of Technology, Tehran, Iran
3. MSc of Information Technology, Faculty of Industrial Engineering, K.N Toosi University of Technology, Tehran, Iran
4. Head of Sarem hospital clinic, Tehran, Iran.

