

## آشکارسازی و تشخیص بیماری دیابت با استفاده از خوشه بندی و تکنیک‌های داده‌کاوی

فاطمه کردی<sup>۱\*</sup>، ابوالفضل اسفندی<sup>۲</sup>، فرزاد همتی<sup>۳</sup>

- ۱- گروه مهندسی کامپیوتر، واحد بروجرد، دانشگاه آزاد اسلامی، بروجرد، ایران
- ۲- دانشگاه آزاد اسلامی واحد بروجرد، باشگاه پژوهشگران جوان و نخبگان بروجرد، ایران
- ۳- گروه مهندسی کامپیوتر، واحد بروجرد، دانشگاه آزاد اسلامی، بروجرد، ایران

### چکیده

در این پژوهش سعی داریم به بررسی و تشخیص بیماری دیابت با استفاده از تکنیک‌های داده‌کاوی و خوشه‌بندی بپردازیم. به‌طور کلی مشکل عمده‌ای که در رابطه با بیماری دیابت وجود دارد عدم تشخیص به‌موقع و یا به‌طور کلی ضعف در تشخیص این بیماری است. داده‌کاوی یکی از بهترین و دقیق‌ترین راه‌ها در پردازش داده‌ها جهت پیش‌بینی می‌باشد. برای تجزیه و تحلیل اطلاعات در این پروژه از نرم‌افزار (MATLAB\R2016a) و ماشین‌بردار پشتیبانی (SVM) به‌منظور تشخیص بیماری استفاده شده است. تعداد داده‌های مورداستفاده ۷۶۸ داده است. این داده‌ها شامل ۹ ویژگی است که هریک بیان‌کننده مشخصه‌ای پزشکی مربوط به وضعیت فرد سالم و بیمار مبتلا به دیابت می‌باشد. با استفاده از الگوریتم خوشه‌بندی k-means بدون نظارت، داده‌های نويز دار دارای شرایط خاص را از دیتاست حذف خواهیم کرد تا دقت کلاس‌بندی بالا برود. نتایج، نشان‌دهنده صحت بالاتر روش پیشنهادی با دقت ۹۴/۲ درصد در این پژوهش در مقایسه با سایر روش‌های هوشمند مورداستفاده تاکنون جهت تشخیص بیماری دیابت می‌باشد.

کلمات کلیدی: الگوریتم k-means، الگوریتم ماشین بردار پشتیبان، بیماری دیابت، خوشه‌بندی، داده‌کاوی

### ۱. مقدمه

دیابت یا بیماری قند یکی از شناخته‌شده‌ترین بیماری‌ها در بین عموم مردم به شمار می‌آید، به‌طور کلی بیماری دیابت رایج‌ترین بیماری غددی و چهارمین علت مرگ در کشورهای پیشرفته می‌باشد و پیشگیری از آن بدون شک موضوع

\* نویسنده اول: گروه مهندسی کامپیوتر، واحد بروجرد، دانشگاه آزاد اسلامی، بروجرد، ایران Email: [f.kordi89@gmail.com](mailto:f.kordi89@gmail.com)

نویسنده دوم: دانشگاه آزاد اسلامی واحد بروجرد، باشگاه پژوهشگران جوان و نخبگان بروجرد، ایران

Email: [abolfazlesfandi@gmail.com](mailto:abolfazlesfandi@gmail.com)

نویسنده سوم: گروه مهندسی کامپیوتر، واحد بروجرد، دانشگاه آزاد اسلامی، بروجرد، ایران Email: [Hemati.f@gmail.com](mailto:Hemati.f@gmail.com)

حیاتی، درمانی و اقتصادی در قرن بیست و یکم است. دیابت می‌تواند منجر به مشکلات حاد فیزیکی در افراد و تأثیر اقتصادی بزرگی بر سیستم بهداشت و درمان ملی گردد [۱]. مخارج درمانی بیماران دیابتی در سال ۲۰۱۰ برابر با ۱۱/۶ درصد مخارج کل جهان بوده است. در این راستا شناسایی و تشخیص دقیق و به‌موقع این بیماری از اهمیت بسیار بالایی برخوردار است [۳ و ۲].

در حالت ساده می‌توان گفت اهمیت پیش‌بینی دیابت از این لحاظ است که بیمار پس از این آگاهی می‌تواند روش زندگی خود را تغییر داده و از این طریق از وقوع بیماری پیشگیری کند. در طی سال‌های اخیر استفاده از تکنولوژی‌های مدرن در علم پزشکی گسترش یافته است، در این پژوهش نیز با در نظر گرفتن موضوع سعی داریم به بررسی و تشخیص بیماری دیابت با استفاده از تکنیک‌های داده‌کاوی و خوشه‌بندی بپردازیم. بیماری دیابت در حال حاضر یکی از شناخته‌شده‌ترین بیماری‌ها به حساب می‌آید، این بیماری به‌طور کلی یکی از مخرب‌ترین بیماری‌ها در جهان به حساب می‌آید که عمدتاً در کشورهای توسعه‌یافته و در حال توسعه در حال گسترش می‌باشد. طبق گزارشی در ایران از هر ۲۰ نفر یک نفر به این بیماری مبتلا می‌باشد وعده‌ی کثیری از این افراد از بیماری خود اطلاع ندارند [۱].

به‌طور کلی مشکل عمده‌ای که در رابطه با بیماری دیابت وجود دارد عدم تشخیص به‌موقع و یا به‌طور کلی ضعف در تشخیص این بیماری است که این ضعف نیز به دلیل عدم انتخاب الگوی مناسب توسط پزشک و یا عدم استفاده مناسب از الگوهای استاندارد است؛ بنابراین پیاده‌سازی روشی که بتواند هر فرد را در تشخیص صحیح ابتلا یا عدم ابتلا به این بیماری یاری رساند می‌تواند گام مهمی در جهت پیشگیری و کنترل این بیماری به‌خصوص در مراحل ابتدایی آن باشد. دیابت همانند سایر بیماری‌ها عوارض بسیاری را به دنبال دارد، گرفتگی قلبی و عروقی و در نوع پیشرفته آن نابینایی، قطع اعضای بدن و اختلالات فکری از جمله عوارض این بیماری می‌باشد، مشکل عمده‌ای که در حال حاضر در رابطه با این بیماری وجود دارد عدم تشخیص به‌موقع و یا به‌طور کلی ضعف در تشخیص این بیماری می‌باشد [۴].

حال با توجه به گسترش این بیماری استفاده از روش‌های جدید پزشکی در راستای بهره‌گیری در تشخیص بسیار مورد توجه قرار گرفته است که استفاده از تکنیک‌های داده‌کاوی و خوشه‌بندی در این پژوهش در جهت تشخیص این بیماری در دستور کار ما قرار دارد، به‌طور کلی داده‌کاوی می‌تواند ارتباطات و وابستگی‌های جدید و بدیعی را کشف کند که برای پزشکان مفید هستند و با استفاده از نتایج کار می‌توان نتایج خوبی را در راستای کمک به علم پزشکی گرفت [۵].

## ۲. خوشه‌بندی

خوشه‌بندی، وظیفه تقسیم یک گروه ناهم‌جنس را در چندین زیرگروه بر عهده دارد. این فرآیند یک تفاوت اساسی با طبقه‌بندی دارد؛ زیرا در این مدل هیچ‌گونه الگوی آموزشی نداریم. خوشه‌بندی به‌طور خودکار ویژگی‌های متمایزکننده زیرگروه‌ها را تعریف می‌کند و زیرگروه‌ها را سازماندهی می‌نماید و به‌عنوان نوعی قابلیت داده‌کاوی غیرمستقیم مطرح است. اغلب از خوشه‌بندی به‌عنوان اولین گام فرآیندهای داده‌کاوی یاد می‌شوند که قبل از سایر فرآیندها برای شناسایی گروهی از رکوردهای مرتبط باهم که بعداً بتوانند نقطه آغاز تحلیل‌ها باشند بر روی رکوردها اعمال می‌شود [۶].

### ۳. تجزیه و تحلیل خوشه‌بندی

تعدادی از تکنیک‌های خوشه‌بندی نظیر الگوریتم‌های خوشه‌بندی داده‌کاوی محور نظیر:

- K-MEANS -
- Y-MEANS -
- FUZZY C-MEANS -

وجود دارند که با توجه به محدودیت‌هایمان در ارتباط با موضوع و تدوین پژوهش به بررسی از الگوریتم K-MEANS می‌پردازیم. همچنین با تلفیق این روش خوشه‌بندی و تکنیک‌های داده‌کاوی به پردازش داده‌های متناسب با پژوهش و تشخیص بیماری دیابت می‌پردازیم [۷].

#### ۳.۱. خوشه‌بندی K-MEANS

می‌توان گفت این الگوریتم به‌سختی قسمت‌بندی شده و در بسیاری جاها به‌سادگی و با سرعت به کار می‌رود. این الگوریتم فاصله اقلیدسی را با اندازه‌های مشابه استفاده می‌کند. مفهوم خوشه‌بندی سخت این است که یک آیتم در مجموعه داده می‌تواند متعلق به یکی و فقط یک خوشه در یک‌زمان باشد. این الگوریتم یک تحلیل خوشه‌بندی است به‌طوری‌که آیتم‌های گروه‌ها بر اساس مقادیر مشخص در داخل K به خوشه‌هایی که آیتم‌ها در همان خوشه ویژگی‌های مشابه دارند ملحق نمی‌شود. چون آن‌ها در خوشه‌های متفاوت مشخصه‌های متفاوتی دارند. عملکرد اقلیدسی با محاسبه فاصله بین دو آیتم به‌صورت زیر استفاده می‌شود.

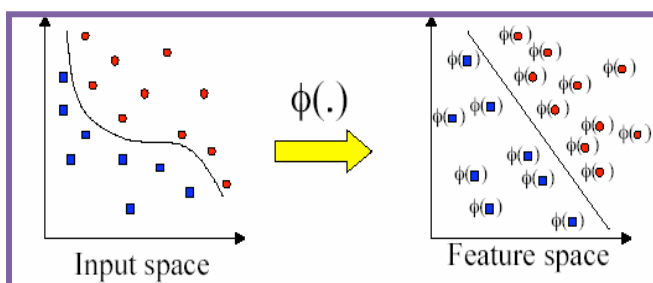
$$d(p,q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2} \quad (1)$$

در اینجا  $p = (p_1, p_2, \dots, p_m)$  و  $q = (q_1, q_2, \dots, q_n)$  دو بردار ورودی با مشخصه‌های  $m$  کمیت هستند. الگوریتم اجرا می‌شود تا مجموعه‌های آموزش داده که ممکن است شامل ترافیک عادی و یا غیرعادی شود بدون اینکه قبلاً برچسب خورده باشد. ایده اصلی این دیدگاه بر اساس فرضیه است که ترافیک عادی و غیرعادی از خوشه‌ها می‌باشد. همچنین داده ممکن است شامل بخش‌های مجزایی شود که کدام آیتم‌های داده از آیتم‌های دیگر در خوشه خیلی متفاوت هستند و کدام متعلق به هیچ خوشه‌ای نیست. بخش مجزا با مقایسه شعاع آیتم‌های داده درک می‌شود به‌طوری‌که اگر شعاع یک آیتم داده بزرگ‌تر از یک آستانه ارائه‌شده باشد. پس به‌صورت یک بخش مجزا دیده می‌شود؛ اما این موضوع یعنی پردازش خوشه‌بندی k-means را که تعداد بخش مجزا و کوچک است را بر هم نمی‌زند [۷].

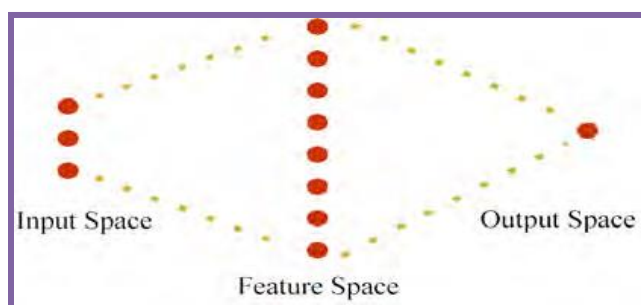
### ۱.۱.۳. روش کار الگوریتم خوشه‌بندی k-means

- تعریف تعداد خوشه‌های  $k$  برای مثال اگر  $k=2$  باشد. ما در آموزش دو خوشه متفاوت، ترافیک عادی و غیرعادی را فرض می‌کنیم.
- مقداره‌ی اولیه مراکز ثقل خوشه  $k$ ، این کار با انتخاب اتفاقی آیتم‌های داده  $k$  از مجموعه داده انجام می‌شود.
- محاسبه فاصله هر آیتم تا مراکز ثقل همه‌ی خوشه‌ها با به‌کارگیری فاصله اقلیدسی متریک که برای یافتن شباهت بین آیتم‌ها در مجموعه داده با کار می‌رود.
- اختصاص هر آیتم با نزدیک‌ترین مرکز ثقل خوشه، در این روش همه‌ی آیتم‌ها به خوشه‌های مختلف اختصاص خواهند یافت، به طوری که هر خوشه آیتم‌هایی با ویژگی‌های مشابه خواهد داشت.
- پس از اختصاص یافتن همه‌ی آیتم‌ها به خوشه‌های مختلف، میانگین خوشه‌های تغییر یافته مجدداً محاسبه می‌شود. میانگین اخیراً محاسبه شده به عنوان مرکز ثقل جدید اختصاص می‌یابد.
- تکرار مرحله ۳ تا زمانی که مرکز ثقل خوشه تغییر نکند.
- برچسب زدن خوشه‌ها به صورت عادی و غیرعادی، به تعداد آیتم‌های داده هر خوشه بستگی دارد.
- می‌توان گفت یک مشکل اساسی در این روش خوشه‌بندی وجود دارد، که این مشکل تعیین بخش اولیه و تعداد مناسب خوشه‌های  $K$  می‌باشد. همچنین گاهی اوقات منجر به هم تراز می‌شود، که میانگین کدام پردازش خوشه بندی ممکن است با چند خوشه خالی به پایان برسد [۷].

### ۴. ماشین بردار پشتیبان



شکل ۱: نگاهت داده‌ها از فضای ورودی به فضای ویژگی



شکل ۲: نمونه‌ای از نوع تبدیل SVM

در این پژوهش از ماشین بردار پشتیبانی SVM (Support Vector Machines) به منظور تشخیص بیماری استفاده شده است. SVM یک الگوریتم یادگیری است، که می‌تواند در کاربردهایی نظیر جداسازی و طبقه‌بندی داده‌ها مورد استفاده قرار گیرد. این الگوریتم پس از طی دو مرحله آموزش می‌تواند داده‌های ورودی را به دسته تقسیم کند و به همین جهت از سادگی منحصربه‌فردی نسبت به سیستم‌های هوشمندی هم چون شبکه‌های عصبی برخوردار است. البته به علت اینکه SVM در جداسازی داده به بیش از دو گروه ناتوان است، کمتر در مواردی که تعداد گروه‌های خروجی بیش از دو گروه باشد استفاده می‌شود.

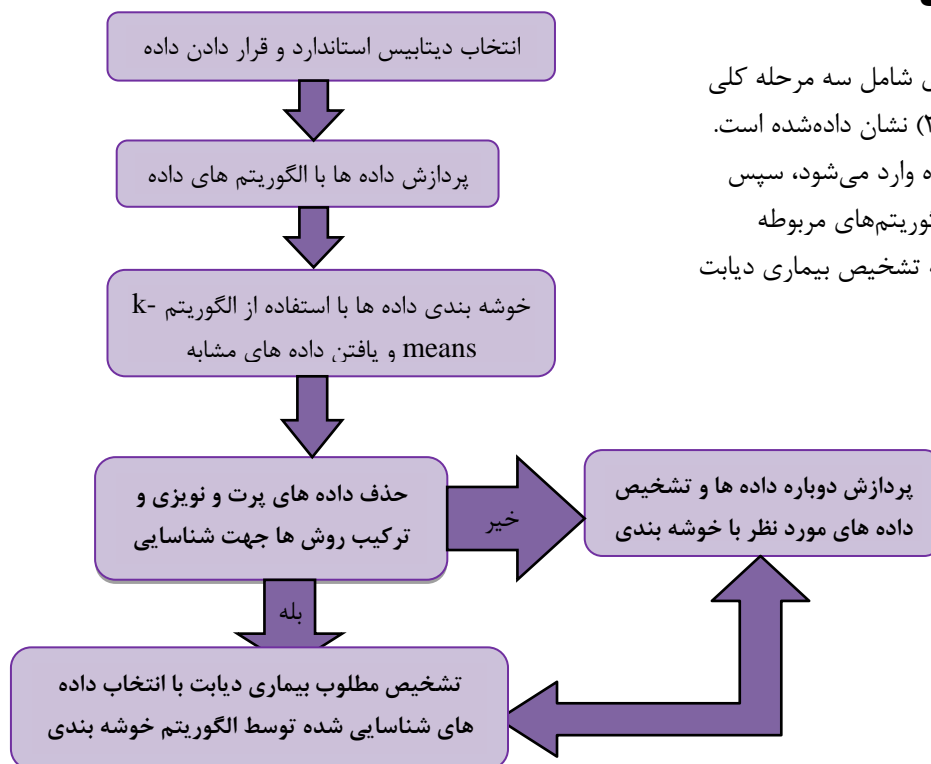
در صورتی که داده در فضای داده‌شده توسط یک خط قابل جداسازی نباشد، الگوریتم با نگاشت ویژگی‌ها به وسیله تابع به ابعاد بالاتر سعی می‌کند توسط یک فرا صفحه با دقت بالاتری داده‌ها را تفکیک نماید. حداکثر حاشیه ایده اولیه SVM خطی بوده است که در آن مطابق شکل (۱) و (۲) طبقه‌بندی داده‌های توسط کرنل خطی صورت می‌گیرد. روش طبقه‌بندی ماشین بردار پشتیبان به این صورت می‌باشد، که با فرض این که مجموعه داده‌ها را به صورت مجموعه زیر نمایش دهیم:

$$[(X_1, C_1), (X_2, C_2), \dots, (X_n, C_n)] \quad (5)$$

$X_i$  یک نمونه از داده‌ها را نشان می‌دهد.  $C_i$  کلاس یا طبقه مربوط به نمونه  $X_i$  را مشخص می‌کند. به عنوان مثال  $C_i$  می‌تواند نشان‌دهنده افرادی که دیابت دارند یا دیابت ندارند، باشد. و آن‌ها را به دو طبقه  $C_i = \{-1, 1\}$  تفکیک کنیم. هر  $X_i$  یک بردار  $p$  بعدی از اعداد حقیقی است. روش‌های طبقه‌بندی خطی، سعی دارند که با ساختن یک ابر سطح (که عبارت است از یک معادله خطی)، داده‌ها را از هم تفکیک کنند. روش طبقه‌بندی ماشین بردار پشتیبان که یکی از روش‌های طبقه‌بندی خطی است، بهترین ابر سطحی را پیدا می‌کند، که با حداکثر فاصله (maximum margin)، داده‌های مربوط به دو طبقه را از هم تفکیک کند [۱۰ و ۱۱].

## ۵. الگوریتم پیشنهادی

الگوریتم پیشنهادی شامل سه مرحله کلی می‌باشد که در شکل (۳) نشان داده شده است. در مرحله اول پایگاه داده وارد می‌شود، سپس به اجرای تکنیک‌ها و الگوریتم‌های مربوطه می‌پردازیم و در نهایت به تشخیص بیماری دیابت دست می‌یابیم.



شکل ۳: الگوریتم پیشنهادی

## ۶. پایگاه داده

هدف این پژوهش تشخیص بیماری دیابت می‌باشد که برای محقق شدن این هدف از پایگاه داده UCİ استفاده کرده‌ایم. این دیتاست شامل اطلاعات زنان با ۹ ویژگی و ۷۶۸ سطر می‌باشد که توسط سازمان بهداشت جهانی معیارهای تشخیص دیابت گردآوری شده و شرح داده‌ها به صورت زیر است:

- ۱- تعداد دفعات بارداری<sup>۱</sup>
- ۲- غلظت گلوکز پلاسما ۲ ساعت در آزمایش تحمل گلوکز خوراکی خون<sup>۲</sup>
- ۳- فشارخون دیاستولیک (mmHg)
- ۴- ضخامت پوست ماهیچه<sup>۳</sup> (mm)
- ۵- انسولین سرم ۲ ساعته<sup>۴</sup> (UU / میلی لیتر)
- ۶- شاخص توده (جرم) بدن<sup>۵</sup> (وزن kg / ((ارتفاع در متر) ^ ۲))
- ۷- داشتن سابقه دیابت<sup>۶</sup>
- ۸- سن<sup>۷</sup>
- ۹- نتیجه آزمایش<sup>۸</sup>

## ۷. ابزارهای پیاده‌سازی

در این پژوهش جهت تشخیص دیابت از روش خوشه‌بندی استفاده شده است و سعی خواهیم کرد تا عمل خوشه‌بندی را با فیلدهای ورودی انجام بدهیم و داده‌های نويز دار را حذف کنیم. در ادامه کار با استفاده از الگوریتم بردار پشتیبان، عمل کلاس‌بندی را نیز انجام خواهیم داد. حال برای خوشه‌بندی تمام ورودی‌ها را با ویژگی‌های ۱ تا ۸ در نظر خواهیم گرفت. بعد از خوشه‌بندی دیتاست به هشت خوشه، در صورتی که در خوشه مورد نظر داده بر چسب‌دار هر کدام از داده‌ها کمتر باشد آن داده بر چسب‌دار را از خوشه مورد نظر حذف خواهیم کرد، بعد از حذف داده‌های نويز دار، دقت کلاس‌بندی داده‌های بدون نويز افزایش قابل توجهی خواهد داشت. حال برای تست و آموزش داده‌ها با بردار پشتیبان، ¼ دیتاست را برای تست در نظر گرفتیم و از بقیه دیتاست برای آموزش شبکه بردار پشتیبان استفاده می‌کنیم که برای پیاده‌سازی الگوریتم پیشنهادی از شبیه‌ساز متلب (MATLAB\R2016a) استفاده شده است.

<sup>1</sup> Numberpregnant

<sup>2</sup> Plasma glucose

<sup>3</sup> skin fold thickness

<sup>4</sup> serum insulin

<sup>5</sup> Body mass

<sup>6</sup> Diabetes pedigree

<sup>7</sup> Age

<sup>8</sup> Class variable

## ۸. ارزیابی الگوریتم‌های دسته‌بندی

در این بخش معیارهای ارزیابی برای دسته‌بندی مورد بررسی قرار خواهند گرفت. ارزیابی یک مدل دسته‌بندی می‌تواند براساس نمونه‌های آموزشی و آزمایشی صورت گیرد. برای ارزیابی باید برچسبی که مدل دسته‌بندی به نمونه‌ها نسبت داده شده است با برچسبی که نمونه‌ها متعلق به آن است مقایسه شود. وقوع حالات مختلف برای دسته‌ها با توجه به مجموعه داده‌های ورودی برای دسته‌بندی با مقادیر TP, FP, TN, FN برای دودسته مثبت و منفی در جدول (۱) نشان داده شده است.

✓ **TP:** مخفف True Positive است و بیانگر تعداد نمونه‌هایی است جدول ۱: حالات هزینه‌ی دسته‌ها

تخمین زده شده		
مثبت		منفی
مثبت	TP	FN
منفی	FP	TN

که دسته واقعی آن‌ها مثبت بوده و الگوریتم دسته‌بندی نیز آن‌ها را به درستی به دسته مثبت نسبت داده است.

✓ **FP:** مخفف False Positive است و بیانگر تعداد نمونه‌هایی است که دسته واقعی آن‌ها منفی بوده و الگوریتم دسته‌بندی آن‌ها را به اشتباه به دسته مثبت نسبت داده است.

✓ **TN:** مخفف True Negative است و بیانگر تعداد نمونه‌هایی است که دسته واقعی آن‌ها منفی بوده و الگوریتم دسته‌بندی نیز آن‌ها را به درستی به دسته منفی نسبت داده است.

✓ **FN:** مخفف False Negative است و بیانگر تعداد نمونه‌هایی است که دسته واقعی آن‌ها مثبت بوده و الگوریتم دسته بندی آن‌ها را به اشتباه به دسته منفی نسبت داده است.

با توجه به پارامترهای مطرح شده در جدول (۱) معیارهای ارزیابی مختلفی ارائه شده است که از جمله مهم‌ترین آن‌ها می‌توان به معیار درستی، دقت، فراخوانی اشاره کرد [۹ و ۸].

## ۹. معیارهای ارزیابی

- ✓ میزان صحت دسته‌بندی (accuracy)
- ✓ دقت تشخیص کل دیابتی مثبت (recall)
- ✓ دقت تشخیص دیابتی مثبت (precision) [۹ و ۸].

### ۹.۱. معیار Accuracy (دقت، صحت)

مهم‌ترین معیار برای تعیین کارایی یک الگوریتم دسته‌بندی می‌باشد. این معیار، دقت کل یک دسته‌بندی را محاسبه می‌کند. این معیار نشان‌دهنده این موضوع است که چند درصد از کل مجموعه داده‌ها به درستی دسته‌بندی شده است.

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}} \times 100\% \quad (2)$$

دو مقدار TP و TN مهم‌ترین مقادیری هستند که باید بیشینه شوند تا کارایی دسته‌بندی به حداکثر برسد [۹و۸].

### ۲.۹ معیار Precision

درصدی از نمونه‌ها را نشان می‌دهد که از میان تمامی نمونه‌ها که توسط دسته‌بندی به آن دسته نسبت داده شده‌اند، درست دسته‌بندی شده‌اند. به عبارتی دقت دسته‌بندی دسته  $i$  را با توجه به کل مواردی نشان می‌دهد که برچسب  $i$  برای نمونه مورد بررسی توسط دسته‌بندی پیشنهاد شده است. نحوه‌ی محاسبه‌ی این معیار نشان داده شده است. اندیس  $i$  در این پارامتر به این مفهوم است که پارامترها باید برای هر دسته  $i$  محاسبه شوند.

$$\text{Precision} = \frac{\text{TP}_i}{\text{TP}_i + \text{FP}_i} \quad (3)$$

### ۳.۹ معیار Recall

برای یک دسته‌بندی، درصد نمونه‌هایی را نشان می‌دهد که از میان تمام نمونه‌ها متعلق به آن دسته، به درستی، دسته‌بندی شده است. به عبارتی دقت دسته‌بندی دسته  $i$  را با توجه به کل نمونه‌های با برچسب  $i$  نشان می‌دهد.

$$\text{Recall}_i = \frac{\text{TP}_i}{\text{TP}_i + \text{FN}_i} \quad (4)$$

نکته قابل توجه این است که معیار Recall کارایی دسته‌بندی را با توجه به تعداد رخداد دسته  $i$  نشان می‌دهد درحالی که معیار Precision اساساً مبتنی بر دقت پیش‌بینی دسته می‌باشد و بیانگر آن است که به چه میزان می‌توانیم به خروجی دسته‌بندی اعتماد کنیم [۹و۸].

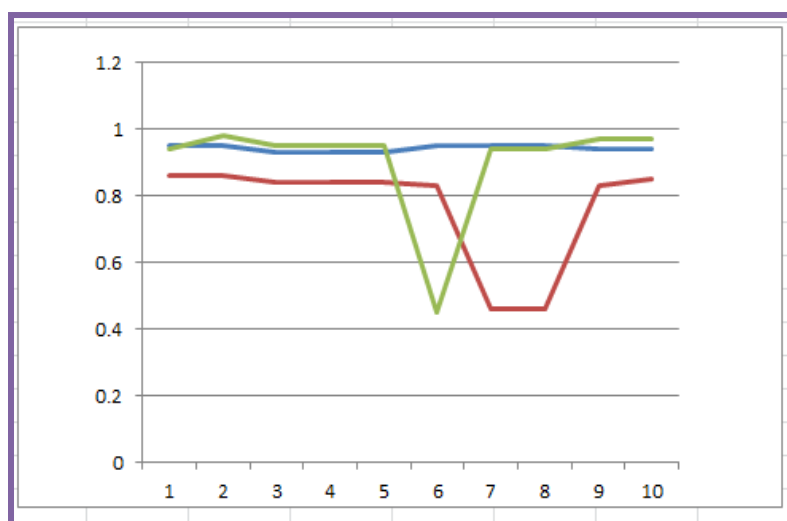


مقدار دقت حساسیت تشخیص با تغییر تکرار الگوریتم را در پژوهش انجام شده به صورت زیر در جدول (۲) می‌باشد.

جدول ۲: دفعات اجرای نمودار به همراه مقادیر سه معیار موردبررسی

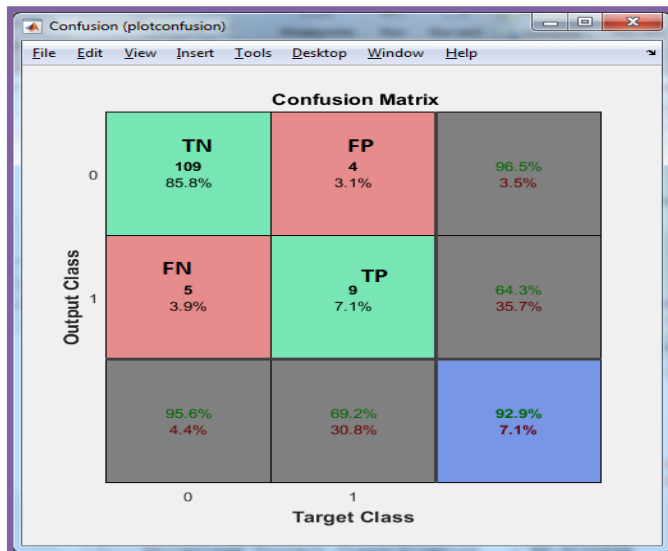
معیار ارزیابی	رنگ قابل نمایش	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰
Accuracy	آبی	۰/۹۴	۰/۹۵	۰/۹۵	۰/۹۵	۰/۹۳	۰/۹۳	۰/۹۳	۰/۹۵	۰/۹۵	۰/۹۴
Recall	قرمز	۰/۸۳	۰/۴۶	۰/۴۶	۰/۸۳	۰/۸۴	۰/۸۴	۰/۸۴	۰/۸۶	۰/۸۶	۰/۸۵
Precision	سبز	۰/۹۷	۰/۹۴	۰/۹۴	۰/۴۵	۰/۹۵	۰/۹۵	۰/۹۵	۰/۹۸	۰/۹۴	۰/۹۷

در شکل (۴) مشاهده می‌کنید که دقت الگوریتم پیشنهادی نرخ ثابتی را دارد اما دو معیار دیگر با افت شاخص مواجه شده است، البته با میانگین‌گیری نرخ ثابتی دارند.



شکل ۴: مقایسه سه معیار با دفعات تکرار ۱۰ بار

## ۱۰. ماتریس آشوب



فصل مشترک تعاریفی که برای مفهوم آشوب ارائه شده است، تأکید بر این نکته است که آشوب دانش بررسی رفتار سیستم‌هایی است که اگرچه ورودی آن‌ها قابل‌تعیین و اندازه‌گیری است، اما خروجی این سیستم‌ها ظاهری تصادفی دارد. شاید به همین دلیل بود که استوارت ریاضیدان برجسته این موضوع را مفهومی احتمالاتی می‌دانست، اما چیزی نگذشت که وی تعریف خود را اصلاح کرد و به تعریفی رسید که تقریباً مورد تأیید عمومی قرار دارد. بر اساس این تعریف، آشوب به توانایی یک الگو و مدل ساده گفته می‌شود که اگرچه خود این الگو نشانی از پدیده‌های تصادفی در خود ندارد، اما می‌تواند منجر به ظهور رفتارهای بسیار در محیط شود. شکل (۵) ماتریس آشوب را نمایش می‌دهد.

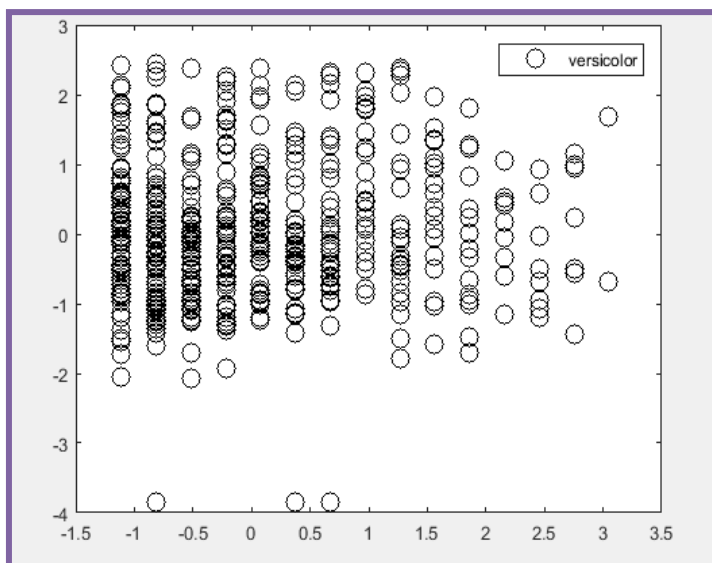
شکل ۵: عملکرد ماتریس آشوب بر روی پروژه

## ۱۱. روش‌های پیاده‌سازی در این پژوهش

- پیاده‌سازی بدون الگوریتم k-means
- پیاده‌سازی با الگوریتم k-means
- کاهش ویژگی‌ها

### ۱۱.۱. پیاده‌سازی بدون الگوریتم k-means

داده‌های مورد آزمایش و تست را تعیین می‌کنیم. نهایتاً دقت صحت کلاس‌بندی را اندازه‌گیری می‌کنیم. بردارهای پشتیبان را نیز در تصویر زیر مشاهده می‌کنید. همچنین با بررسی ماتریس آشوب می‌توان تعداد کلاس‌بندی اشتباه را مشاهده کرد. ۰/۴۱۱۵، دقت کلاس‌بندی است که بر اساس تعداد درستی برنامه تشخیص داده شده است. این از روش ¼ داده‌ها ۳۴ مورد است یعنی صحت تشخیص درست مثبت ۷۲/۹۱ درصد می‌باشد. در شکل (۶) خروجی بردار پشتیبان نمایش داده شده است.



در خروجی بردار پشتیبان، محور X فضای چندبعدی است که شامل چندین ویژگی خواهد بود و محور Y خروجی با برچسب داده ها خواهد بود. در واقع به تعداد ویژگی ها برچسب تولید می کند، دایره های خالی متعلق به خوشه بندی ۱، دایره های توپر هم برای خوشه بندی ۰ است.

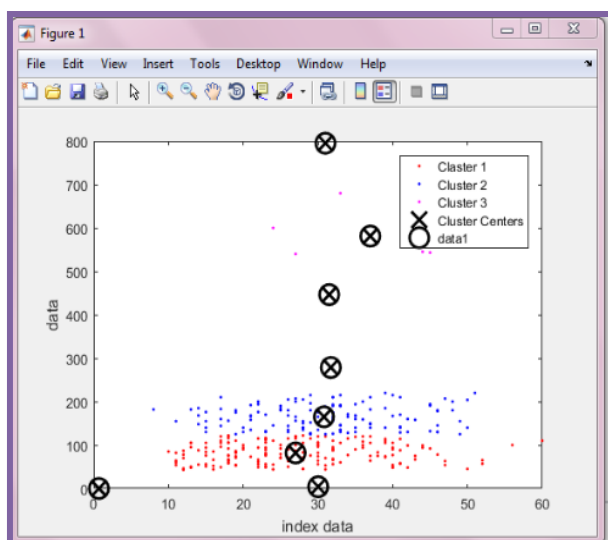
شکل ۶: خروجی بردار پشتیبان بدون الگوریتم k-means

دقت کلاس بندی در حالت کلی به صورت زیر می باشد:

Percentage Correct Classification: 72. 916667%

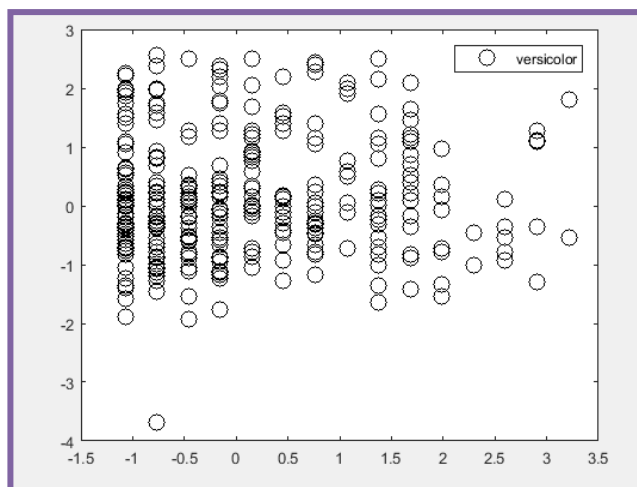
Percentage Incorrect Classification: 27. 083333%

### ۲.۱۱. پیاده سازی با الگوریتم k-means



در این روش از کار با استفاده از الگوریتم خوشه بندی k-means داده ها را به ۸ خوشه تقسیم کرده همه خوشه ها را برای حذف داده های نویز دار بررسی می کنیم. در هر خوشه با بررسی داده های برچسب دار متعلق به یک کلاس در صورتی که از دیگر داده ها کمتر باشد، به عنوان نویز شناخته و از دیتاست حذف می کنیم. بعد از آماده سازی دیتاست داده ها، عمل کلاس بندی را با استفاده از بردارهای پشتیبان انجام داده و دقت کلاس بندی را اندازه گیری می کنیم  $\frac{1}{4}$  دیتاست را به عنوان تست در نظر می گیریم.

شکل ۷: خروجی خوشه بندی



شکل ۸: خروجی بردار پشتیبان با الگوریتم k-means

بعد از هر بار اجرا حدود ۳۰۰ سطر از داده نویز دار را حذف کرده، صحت دیتاست را بالا می‌بریم. بعد از کلاس‌بندی میزان خطا به صورت قابل توجهی کاهش یافته است. تابع K-means مجموعه را به  $k$  خوشه‌ی کاملاً مجزا تبدیل می‌نماید و اندیس هر خوشه را به داده‌های مجموعه نسبت می‌دهد. نمونه‌ها معمولاً به وسیله بردارهایی (نقاطی) در یک فضای چندبعدی نمایش داده می‌شوند که هر بعد یک صفت (متغیر) مشخص از نمونه‌ها را نشان می‌دهد. مانند سن، فشارخون و غیره اگر فرض کنیم در شکل (۷) و (۸) خروجی خوشه‌بندی و بردار پشتیبان نمایش داده شده است.

اگر  $n$  نمونه داشته باشیم که هر کدام دارای  $d$  صفت باشند، پس یک ماتریس نمونه  $n \times d$  را می‌توانیم تشکیل دهیم که هر سطر این ماتریس یک نمونه و هر ستون یک صفت مربوط به نمونه‌ها می‌باشد. داده‌ها با استفاده از روش min-max نرمال‌سازی می‌شوند، نرمال‌سازی تغییر مقیاس داده‌ها به گونه‌ای است که آن‌ها را به یک فاصله کوچک و معین نگاشت می‌کند. مطابق شکل نه تنها مراکز خوشه‌ها ( $\otimes$ ) مشخص شده است، بلکه هر خوشه بارنگی خاص نمایش داده شده و محور  $X$  در این شکل سن افراد و محور  $Y$  تعداد داده‌هایی که در این خوشه قرار گرفته‌اند را مشخص می‌کند. نمایی از مراکز خوشه‌های هشت‌گانه از چهل و یک صدم به ۲۱ صدم کاهش یافته است.  $0.2632$  دقت کلاس‌بندی است که براساس تعداد درستی برنامه تشخیص داده شده است و میزان خطای کلاس‌بندی که کاهش یافته است، با بررسی ماتریس آشوب صحت کلاس‌بندی ۹۰ درصد است.

Percentage Correct Classification: 93. 181818%

Percentage Incorrect Classification: 6. 818182%

### ۳.۱۱. کاهش ویژگی‌ها

حال با تست حالت‌های مختلف، بهترین حالت ورودی را برای کاهش خطا و افزایش دقت انتخاب کردیم، که از ویژگی‌های ۸ تا ۴ به ترتیب استفاده شده است:

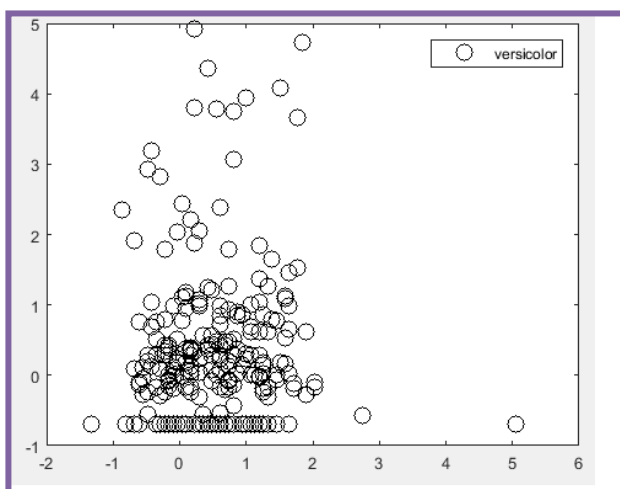
۴- skin fold thickness (ضخامت پوست ماهیچه (mm))

۵- serum insulin (انسولین سرم ۲ ساعته (UU / میلی‌لیتر))

۶- Body mass (شاخص توده بدن (وزن kg / (ارتفاع در متر)  $^2$ ))

۷- Diabetes pedigree (داشتن سابقه دیابت)

۸- Age (سن)



به‌عنوان ورودی در نظر گرفتیم، که از الگوریتم k-means برای حذف داده‌های نویز دار استفاده می‌کنیم. صحت تشخیص کلاس‌بندی افزایش قابل‌ملاحظه‌ای داشت، که در حالت بدون حذف نویز تأثیر چندانی داشت. در روش اولی ذکر شده صحت تشخیص کلاس‌بندی ۹۳/۷ درصد افزایش یافت. میزان خطا نیز به ۱۷ درصد رسید. در شکل (۹) نیز تفکیک داده‌ها با بردار پشتیبان قابل‌مشاهده است.

شکل ۹: خروجی بردار پشتیبان با الگوریتم k-means و کاهش ویژگی‌ها

Percentage Correct Classification: 93. 700787%

Percentage Incorrect Classification: 6. 299213%

## ۱۲. مقایسه روش پیشنهادی با سایر روش‌ها

در این جدول (۳) به‌صورت خلاصه روش‌های بررسی‌شده در خصوص بیماری دیابت با الگوریتم‌های داده‌کاوی و درصد درستی آن‌ها در مقایسه با روش پیشنهادی می‌باشد.

جدول ۳: روش‌های بررسی‌شده بیماری دیابت با داده‌کاوی

محقق	روش	Accuracy
پلیت و گون <sup>۱</sup> (۲۰۰۸)	LS – SVM (10- fold CV)	۷۸/۲۱
ایسا و مامات <sup>۲</sup> (۲۰۱۱)	GDA-LS – SVM (10- fold CV)	۸۲/۰۵
آیبینو و همکاران (۲۰۱۱)	Clustered – HMLP	۸۰/۵۹
پلیت و گونز (۲۰۰۷)	Combining PCA and ANFIS	۸۱/۲۸
قهرمانی و اله وردی (۲۰۰۸)	Hybrid system (ANN and FNN )	۸۴/۲
پتیل (۲۰۱۰)	Hybrid Prediction Model (HPM)	۹۲/۳۸
فاطمه کردی (۲۰۱۷)	svm with kmeans	۹۴/۲

<sup>۱</sup>Polat and Gunes

<sup>۲</sup> Isa&Mamat

### ۱۳. نتیجه‌گیری

در این پژوهش ما یک روش با خطای کمتر نسبت به سایر روش‌های تشخیص دیابت توسعه دادیم. این دیتاست شامل داده‌های نویز دار و داده‌های تکراری است. در این روش ما با استفاده از الگوریتم خوشه‌بندی k-means برای تشخیص داده‌های نویز دار استفاده کردیم. بعد از پاک‌سازی داده‌ها، داده‌های بدون نویز داریم. داده‌ها بیشترین تأثیر را در صحت کلاس‌بندی با دقت بالا را دارند. در ادامه کار که داده‌ها را به وسیله SVM کلاس‌بندی کردیم، نسبت به سایر روش‌های قبلی کار شده از صحت و دقت بالایی برخوردار است. همچنین با انتخاب ویژگی‌های مهم تأثیر زیادی در صحت کلاس‌بندی داده‌ها دارند. میزان صحت تشخیص داده‌ها به مرز ۹۴/۲ درصد رسید.

### ۱۴. منابع و مأخذ

۱. مرجب، سونا. و فاضلی، امین. (۱۳۹۵)، " استخراج ویژگی و تشخیص هوشمند بیماری رتینوپاتی دیابتی در تصاویر شبکیه"، سومین کنفرانس ملی توسعه علوم مهندسی، مازندران، تنکابن، موسسه آموزش عالی آیندگان.
2. Bose, I.(2006), "*Data Mining in Diabetes Diagnosis and Detection*", Encyclopedia of Data Warehousing and Mining: IZ, p. 257.
3. Huang ,Y. P. and Black, N. and Harper, R.(2007), "*Feature selection and classification model construction on type 2 diabetic patients' data*", Artificial intelligence in medicine, vol. 41, pp. 251-262.
4. Sankar, K. Pal. and Pabitra, Mitra.(2009), "*Pattern Recognition Algorithms for Data Mining*", (Book), Chapman & Hall/CRC.
5. Acharya, U. R. and Faust, O. and Sree, S.V. and Ghista, D.N. and Dua, S. and Joseph, P. and Tamura, T.(2013), "*An integrated diabetic index using heart rate variability signal features for diagnosis of diabetes*", Computer methods in biomechanics and biomedical engineering, 16(2), 222-234.
6. Snehal, A. Mulay. And Devale, P. R. and Garje, G.V. (2010), "*Intrusion Detection System using Support Vector Machine and Decision tree*" international journal computer application (0975 –8887) Volume 3 –No.3.
7. Sonukumari, A.(2013) , "*A data mining approach for the Diagnosis of Diabetes Mellitus*". 7<sup>th</sup> International Conference on Intelligent Systems and Control, Coimbatore , TamilNadu, India.
8. Domingos, P.(1999), "*MetaCost: a general method for making classifiers cost-sensitive*," in Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 155-164.



9. Engelgau, Michael M. and Narayan, K.M. Venkat . and Herman William H.(2000), "*Screening for type 2 diabetes*", Diabetes care, vol. 23, pp. 1563-1580.
10. Han,J. and Kamber,M. and Pei,J.(2006), " *Data mining: concepts and techniques*" Morgan kaufmann.
11. Xue,H. and Yang,Q. and Chen,S .(1999),"*SVM: Support Vector Machines*"