



دانشگاه کاشان  
University of Kashan

مجله محاسبات نرم

## SOFT COMPUTING JOURNAL

تارنمای مجله: [scj.kashanu.ac.ir](http://scj.kashanu.ac.ir)



### تشخیص بیماری دیابت با استفاده از مدل رای گیری نرم

سکینه اسدی امیری<sup>1\*</sup>، استادیار، هانا یوسف پور<sup>1</sup>، دانشجوی کارشناسی ارشد، سعیده محمدپور<sup>1</sup>، دانشجوی کارشناسی ارشد، دانشجوی کارشناسی ارشد<sup>1</sup> گروه مهندسی کامپیوتر، دانشکده مهندسی و فناوری، دانشگاه مازندران، بابلسر، ایران.

#### اطلاعات مقاله

#### چکیده

تاریخچه مقاله:

دریافت 15 اسفند ماه 1402

پذیرش 9 آبان ماه 1403

کلمات کلیدی:

یادگیری ماشین

انتخاب ویژگی

داده کاوی

دیابت

پیش بینی

دیابت یکی از عوامل مهم مرگ و میر در سراسر جهان است و تاثیرات آن بر بیماری های کلیوی و قلبی و از دست دادن بینایی قابل توجه است. پیش بینی دیابت یک حوزه تحقیقاتی مهم است که می تواند به بهبود درمان بیماری کمک کند. در این مقاله، روش جدیدی برای تشخیص بیماری دیابت پیشنهاد شده است. روش پیشنهادی روی مجموعه داده دیابت اعمال شده است، ابتدا در مرحله پیش پردازش، شناسایی داده های پرت و حذف آنها، جایگزین نمودن مقادیر گمشده و نرمال سازی داده ها انجام می شود. پس از پیش پردازش داده ها با استفاده از الگوریتم لاسو، ویژگی های مهم انتخاب می شوند. سپس با استفاده از سه طبقه بند K-نزدیکترین همسایه، تقویت گرادبان شدید و کت بوست، نمونه ها به دو کلاس بیماران دیابتی و سالم طبقه بندی می شوند. در پایان برای بهبود روش پیشنهادی از الگوریتم رای گیری نرم برای ادغام سه طبقه بند استفاده شده است. مدل پیشنهادی در این پژوهش با استفاده از معیارهای ارزیابی دقت، صحت و پوشش مورد ارزیابی قرار گرفت. این مدل به دقت 94/4٪، صحت 96/5٪ و پوشش 92/7٪ دست یافت. نتایج حاکی از آن هستند که مدل پیشنهادی با افزایش دقت در تشخیص بیماری دیابت نسبت به سایر، عملکرد بهتری داشته است. بنابراین، با استفاده از این مدل، می توان افرادی که در معرض خطر ابتلا به دیابت هستند را با دقت بیشتری شناسایی کرد و اقدامات پیشگیرانه ای را برای کنترل بیماری دیابت انجام داد.

© 1403 نویسندگان. مقاله با دسترسی آزاد تحت مجوز CC-BY

#### 1. مقدمه

انسولین انجام می شود که در پانکراس تولید می شود. این بیماری به دو نوع تقسیم می شود. نوع اول، یا دیابت نوع 1، به دلیل عدم ترشح انسولین رخ می دهد و به طور معمول در سنین کودکی و نوجوانی بروز می کند. نوع دوم، یا دیابت نوع 2، به دلیل مقاومت بدن به انسولین که ممکن است به دلیل چاقی، رخ دهد. توده بدنی بالا یکی از عوامل مهمی است که می تواند بر سلامت انسان تاثیر بگذارد و ممکن است منجر به مشکلاتی مانند دیابت شود. بارداری یک دوره حساس و مهم در زندگی زنان است و نیازمند مراقبت ویژه است. افزایش وزن در طول بارداری طبیعی

دیابت، یک بیماری مزمن است که در آن بدن توانایی کاهش قند خون را از دست می دهد. قند خون انرژی لازم برای فعالیت های بدن را فراهم می کند و اگر سطح آن افزایش یابد، بدن باید برای کاهش آن اقدام کند. کاهش قند خون توسط یک هورمون به نام

✦ نوع مقاله: پژوهشی

\* نویسنده مسئول

پست(های) الکترونیک: s.asadi@umz.ac.ir (اسدی امیری)

h.yousefpour14@umail.umz.ac.ir (یوسف پور)

s.mohammadpour11@umail.umz.ac.ir (محمدپور)

می باشد. در این پژوهش، با توجه به اهمیت این مساله، سعی شده است تا دقت در شناسایی افراد در معرض خطر بیماری دیابت افزایش داده شود. در روش پیشنهادی ابتدا داده های پرت شناسایی و حذف می شوند. در مرحله بعد از الگوریتم لاسو برای انتخاب ویژگی های مهم و حذف ویژگی های غیر ضروری استفاده شده است. پس از آن، از سه طبقه بند K-نزدیکترین همسایه<sup>1</sup> (KNN)، تقویت گرادیان شدید<sup>2</sup> (XGBoost) و کت بوست (CatBoost) برای طبقه بندی نمونه ها به دو کلاس دیابتی و غیر دیابتی استفاده شده است. در انتها برای بهبود روش پیشنهادی از مدل رای گیری نرم استفاده شده است. رای گیری نرم<sup>3</sup> یک تکنیک قدرتمند برای ادغام پیش بینی های کسب شده توسط الگوریتم های یادگیری ماشین است.

این پژوهش شامل پنج بخش است. بخش اول، مقدمه ای است که در آن تعریف دیابت، سیر بیماری و اهمیت آن بیان شده است. بخش دوم به بررسی پیشینه ادبیات و کارهای انجام شده در زمینه تشخیص و بهبود دقت در تشخیص دیابت می پردازد. بخش سوم به تعریف روش کار در این پژوهش می پردازد که شامل پیش پردازش، انتخاب ویژگی و مدل پیشنهادی است. بخش چهارم نتایج حاصل از الگوریتم پیشنهادی این پژوهش را بررسی می کند. در نهایت، بخش پنجم به نتیجه گیری می پردازد.

## 2. پیشینه ادبیات

با توجه به اهمیت تشخیص زودهنگام بیماری دیابت، تحقیقات گسترده ای با استفاده از مدل های یادگیری ماشین برای تشخیص و پیشگیری از دیابت انجام شده است. در مرجع [5]، رویکردی نوین با استفاده از تکنیک های داده کاوی برای افزایش دقت پیش بینی و قدرت تعمیم بالای پیش بینی های مدل ارائه شده است. این پژوهش پس از پیش پردازش مجموعه داده با استفاده از الگوریتم K-means بهبود یافته و رگرسیون لجستیک، به قدرت تشخیص خوبی دست یافته است. در مرجع [6]، برای ایجاد یک

است و به طور معمول بین 11 تا 16 کیلوگرم است. اما اگر افزایش وزن بیش از حد باشد، ممکن است منجر به مشکلاتی مانند دیابت شود. اگر دیابت در طول بارداری رخ دهد، آن را دیابت بارداری می نامند که می تواند برای مادر و جنین او خطرناک باشد [1]. دیابت، بیماری مزمنی است که در آن سطح بالای قند خون می تواند خطرناک باشد. این بیماری با علائم کوتاه مدتی مانند تشنگی، خستگی، تحریک پذیری، گرسنگی شدید و تکرر ادرار همراه است. در بلند مدت، می تواند منجر به امراض دیگری مانند آسیب تدریجی به کلیه ها، چشم ها و قلب شود. این تهدیدات می توانند تاثیر مهمی بر سطح بهداشت و زندگی افراد مبتلا به این بیماری داشته باشند.

دیابت به مرور زمان اتفاق می افتد و علائمی از خود نشان می دهد. شناسایی این علائم می تواند به تشخیص دیابت در مراحل اولیه کمک کند. با تشخیص این علائم و افزایش فعالیت بدنی و تغذیه سالم، افراد می توانند سلامت خود را حفظ کنند. با این حال، با وجود اهمیت فراوان تشخیص دیابت در مراحل اولیه، پزشکان هنوز نتوانسته اند روش موثقی برای تشخیص این علائم ارائه دهند. بنابراین، شناسایی مراحل اولیه دیابت یکی از چالش های مهم دنیای امروز است [2]. با ورود الگوریتم های تشخیصی به عرصه پزشکی، انقلابی در فرآیند تشخیص زودهنگام بیماری ها رخ داد. این الگوریتم ها با استفاده از داده های جمع آوری شده از مبتلایان به بیماری مورد نظر، برای شناسایی عوامل موثر در بیماری بکار رفته اند. این روش ها می توانند الگوهای پنهان و اطلاعات مفیدی را از داده های حجیم استخراج کنند و به پزشکان در بهبود تشخیص دیابت کمک کنند [3]. استفاده از داده های سلامت الکترونیکی، داده های سنسوری و داده های بالینی می تواند به پزشکان در کشف الگوهای جدید در تشخیص و پیش بینی دیابت کمک کند. همچنین، این روش ها می توانند در پیش بینی عوارض احتمالی دیابت و کمک به بیماران در مدیریت بیماری مفید باشند. به طور کلی، استفاده از استخراج داده و الگوریتم های یادگیری می تواند به بهبود قابل توجه در تشخیص دیابت منجر شود [4].

دقت در تشخیص دیابت یکی از چالش های علم پزشکی

<sup>1</sup> K-Nearest Neighbors

<sup>2</sup> Extreme Gradient Boosting

<sup>3</sup> Soft Voting

یادگیری K-نزدیک‌ترین همسایه، آدابوست و LightGBM برای پیش‌بینی دیابت استفاده کردند. سپس با استفاده از رای‌گیری نرم بین این سه الگوریتم، دسته مناسب هر نمونه را تعیین کردند. مرجع [15] برای پیش‌بینی بیماری دیابت از رای‌گیری نرم بین سه الگوریتم پرسپترون چند لایه، ماشین بردار پشتیبان و K-نزدیک‌ترین همسایه استفاده کردند. مرجع [16] از الگوریتم‌های رگرسیون لجستیک، ماشین بردار پشتیبان، K-نزدیک‌ترین همسایه، تقویت گرادیان، بیز ساده، و جنگل تصادفی برای پیش‌بینی بیماری دیابت استفاده کردند. سپس از رای‌گیری سخت برای ترکیب سه الگوریتم بیز ساده، جنگل تصادفی و ماشین بردار پشتیبان استفاده کردند. با توجه به اینکه الگوریتم K-نزدیک‌ترین همسایه در اغلب پژوهش‌های انجام شده، الگوریتمی امیدبخش در تشخیص بیماری دیابت بوده است، می‌توان از این الگوریتم و دیگر الگوریتم‌های تجمعی استفاده کرد که قدرت تشخیص خوب خود را در داده‌های پیچیده ثابت کرده‌اند. بنابراین، استفاده از این الگوریتم‌ها امکان ساخت مدل تشخیصی با دقت بالاتری را فراهم می‌آورد.

### 3. مجموعه داده

مجموعه داده‌ای که در این پژوهش استفاده شده، مربوط به دیابت در بین زنان قوم پیمان<sup>2</sup> در هند است. این مجموعه داده توسط موسسه ملی دیابت و بیماری‌های گوارشی و کلیوی آماده شده است [17]. هدف از این مجموعه داده، پیش‌بینی ابتلا به دیابت بر اساس اندازه‌گیری‌های تشخیصی خاص است که در مجموعه داده وجود دارد. مجموعه داده شامل 768 نمونه، هشت ویژگی پیش‌بینی‌کننده و یک ویژگی که نشان‌دهنده کلاس (هدف) است. در این پژوهش، 70 درصد از مقادیر مجموعه داده برای آموزش، 10 درصد برای اعتبارسنجی و 20 درصد باقیمانده برای آزمایش استفاده شده است. جدول (1) نشان‌دهنده خلاصه آماری مجموعه داده دیابت هندی پیمان می‌باشد. متغیرهای پیش‌بینی‌کننده شامل تعداد حاملگی‌ها، گلوکز، فشار خون،

مدل تشخیصی با قدرت تعمیم بالا، از مجموعه داده‌ای از بنگلادش استفاده شده و الگوریتم‌های بیز ساده، رگرسیون لجستیک و جنگل تصادفی مورد بررسی قرار گرفته است. همچنین در مرجع [7]، مدلی با استفاده از روش‌های داده‌کاوی برای پیش‌بینی موارد دیابت ارائه شده است. در نهایت، در این پژوهش با استفاده از چندین روش، از جمله J48<sup>1</sup>، بیز ساده و ماشین بردار پشتیبان، پیش‌بینی برای تشخیص دیابت انجام می‌شود. در مرجع [8]، برای پیش‌بینی وقوع دیابت در اندونزی، از مجموعه داده‌ای استفاده شده است که شامل هفت عامل خطر دیابت، از جمله سن، جنسیت، شاخص توده بدنی، سابقه خانوادگی دیابت، فشار خون، مدت زمان دیابتی و سطح گلوکز خون به عنوان عوامل کلیدی در شناسایی دیابت است. در نهایت، از روش خوشه‌بندی K-means و روش طبقه‌بندی بیز ساده، برای تشخیص بیماری دیابت استفاده شده است. در مرجع [9]، با استفاده از اطلاعات بیماران بیمارستان نیروی دفاع بحری و بررسی داده‌ها با استفاده از روش‌های داده‌کاوی، تلاش شده است تا شروع بیماری دیابت را پیش‌بینی کنند. مرجع [10] یک تجزیه و تحلیل جامع را در راستای پیش‌بینی دیابت با استفاده از یادگیری ماشین، روش‌های داده‌کاوی و ابزارهای مرتبط انجام داده است. این پژوهش با استفاده از سه روش مختلف داده‌کاوی، از جمله بیز ساده، ماشین بردار پشتیبان و درخت تصمیم، به بررسی موضوع پرداخته است. مرجع [11] در راستای پیش‌بینی احتمال بروز بیماری قلبی برای بیماران دیابتی، با استفاده از روش‌های پیشرفته احتمالی و بر اساس دقت پیش‌بینی‌شان، ارزیابی و تحلیل انجام داده است.

مرجع [12] برای شناسایی بهترین طبقه‌بند برای ارزیابی نتایج بالینی، از یک مجموعه استاندارد معیارها و شبکه‌های بیز ساده و توابع پایه شعاعی استفاده کرده است. در مرجع [13]، با مقایسه روش‌های مختلف یادگیری ماشین، مدلی برای بهبود عملکرد پیشنهاد شده است. روش‌های پیش‌بینی مانند K-نزدیک‌ترین همسایه و رگرسیون لجستیک، از جمله الگوریتم‌های مورد استفاده در این پژوهش بوده‌اند. مرجع [14]، از سه الگوریتم

<sup>2</sup> PIMA

<sup>1</sup> Java implementation of the C4. 5 algorithm

ضخامت پوست، انسولین، شاخص توده بدنی، فاکتورهای مشخص کننده دیابت و سن می باشد. این مجموعه داده دارای دو کلاس، یعنی ابتلا و عدم ابتلا به دیابت، است. شکل (1) نشان دهنده ماتریس همبستگی ویژگی ها با یکدیگر است. بیشتر است.

جدول (1): خلاصه آماری مجموعه داده دیابت هندی پیم (منبع: محاسبات محقق).

| سن     | فاکتورهای مشخص کننده دیابت | شاخص توده بدنی | انسولین | ضخامت پوست | فشار خون | گلوکز   | تعداد حاملگی ها | شاخص ها      |
|--------|----------------------------|----------------|---------|------------|----------|---------|-----------------|--------------|
| 768.0  | 768.0                      | 768.0          | 768.0   | 768.0      | 768.0    | 768.0   | 768.0           | شمارش        |
| 33.240 | 0.471                      | 31.992         | 79.799  | 20.536     | 69.105   | 120.894 | 3.845           | میانگین      |
| 11.760 | 0.331                      | 7.88           | 115.2   | 15.95      | 19.35    | 31.97   | 3.36            | انحراف معیار |
| 21.0   | 0.078                      | 0.0            | 0.0     | 0.0        | 0.0      | 0.0     | 0.0             | کمینه        |
| 24.0   | 0.24375                    | 27.3           | 0.0     | 0.0        | 62.0     | 99.0    | 1.0             | %25          |
| 29.0   | 0.3725                     | 32.0           | 30.5    | 23.0       | 72.0     | 117.0   | 3.0             | %50          |
| 41.0   | 0.62625                    | 36.6           | 127.25  | 32.0       | 80.0     | 140.2   | 6.0             | %75          |
| 81.0   | 2.42                       | 67.1           | 846.0   | 99.0       | 122.0    | 199.0   | 17.0            | بیشینه       |



شکل (1): ماتریس همبستگی ویژگی ها در مجموعه داده پیم (منبع: محاسبات محقق).

شده، پراکندگی داده‌ها در مجموعه داده را نشان می‌دهد. بر اساس نمودار پراکندگی که در شکل (2) نشان داده شده است، داده‌ها به صورت جفتی با تجمع شبیه‌ترین خروجی‌ها همراه هستند. بیشترین همبستگی در ویژگی‌ها بین تعداد حاملگی‌ها و سن، ضخامت پوست و شاخص توده بدنی و گلوکز و انسولین مشاهده می‌شود. این همبستگی بر اساس ارقام مثبت نمودار پراکندگی است. همچنین، با توجه به ماتریس همبستگی که رابطه قوی بین هیچ یک از ویژگی‌ها و خروجی را نشان نمی‌دهد، در نتیجه هیچ یک از ویژگی‌ها به تنهایی قادر به پیش‌بینی بیماری دیابت نمی‌باشد.

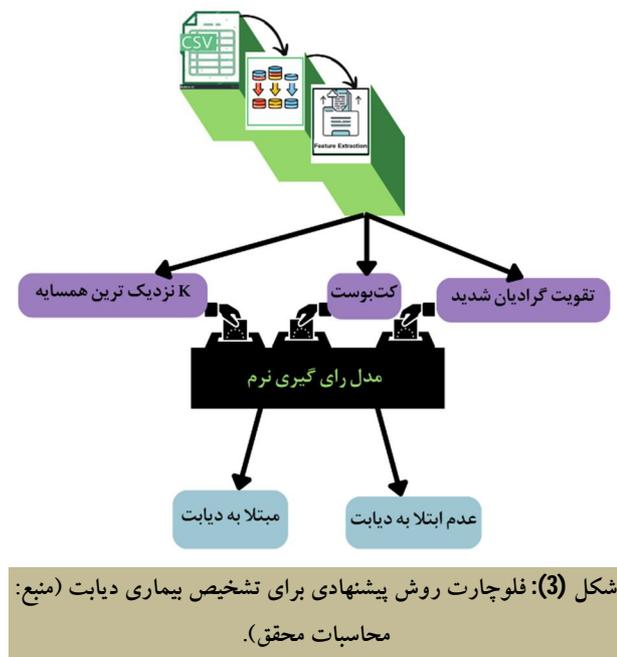
شکل (2) نمودار جفت شده (Pair Plot) ویژگی‌های مجموعه داده بر اساس کلاس را نشان می‌دهد که این ویژگی‌ها از بالا به پایین و از چپ به راست به ترتیب جدول می‌باشند. این نمودار امکان کشف رابطه بین ویژگی‌های موجود در مجموعه داده را فراهم می‌کند. اگر نقاط در نمودار پراکنده باشند، به معنی عدم وجود رابطه آشکاری است. در مقابل، اگر نقاط تقریباً در یک خط مستقیم قرار گیرند، نشان‌دهنده رابطه خطی بین آنها است. در یک نمودار جفت شده، قطرها نمودارهایی هستند که یک ویژگی را در برابر خودش نشان می‌دهند، که این می‌تواند توزیع داده‌ها برای آن ویژگی خاص را نشان دهد. قطر نمودار جفت



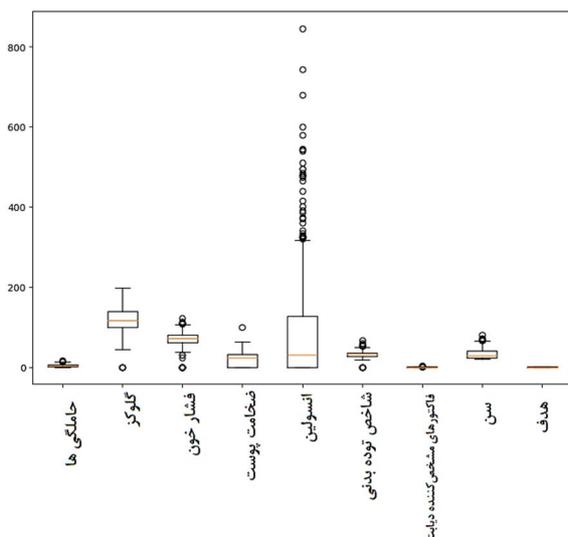
شکل (2): نمودار جفت شده ویژگی‌های مجموعه داده (منبع: محاسبات محقق).

#### 4. روش شناسی

در این مقاله، روش جدیدی برای تشخیص بیماری دیابت معرفی شده است. در مرحله پیش پردازش، مقادیر گمشده با صفر جایگزین می شوند و سپس داده‌ها نرمال سازی<sup>1</sup> می شوند. در مرحله بعدی، با استفاده از روش انتخاب ویژگی لاسو، اهمیت هر یک از ویژگی‌ها محاسبه و ویژگی‌های با اهمیت بیشتر انتخاب می شوند. در نهایت، یک روش با استفاده از رای گیری نرم و سه رای دهنده، یعنی K-نزدیک ترین همسایه، تقویت گرادیان شدید و کت بوست، برای افزایش دقت در تشخیص بیماری دیابت ارائه شده است. شکل (3) فلوجارت روش پیشنهادی در تشخیص بیماری دیابت را نشان می دهد.



داده‌های پرت تفاوت قابل توجهی با بقیه اعضای مجموعه دارند. از آنجا که داده‌های پرت می‌توانند فرآیند تشخیص را برای یادگیرنده دشوار کند، در این پژوهش داده‌های پرت حذف شدند. روش مورد استفاده برای شناسایی داده‌های پرت، روش فاصله بین چارکی<sup>2</sup> (IQR) می‌باشد. هر آنچه که در نمودار جعبه‌ای از چارک دامنه کمینه و بیشینه خارج باشد، داده پرت محسوب می‌شود. در نهایت مقادیر داده‌های پرت با میانگین داده‌ها جایگزین می‌شوند. در این پژوهش، در مرحله پیش پردازش، مقادیر گم شده با مقدار ثابت NaN، با مقدار صفر جایگزین می‌شوند. سپس، داده‌ها نرمال سازی می‌شوند. این عمل، با کاهش دامنه ویژگی‌های مجموعه داده، داده‌ها را به یک دامنه مشخص می‌برد. در این مرحله، دامنه اعداد بین 1 و 1- قرار داده شده است.



شکل (4): نمودار جعبه‌ای نمونه‌های موجود در مجموعه داده (منبع: محاسبات محقق).

#### 1.4. پیش پردازش داده‌ها

قبل از آموزش هر الگوریتم، پیش پردازش داده‌ها ضروری است. پیش پردازش، یکی از مهم‌ترین بخش‌های داده‌کاوی است که به بررسی داده‌ها برای یافتن مقادیر پرت و گم شده و بررسی ویژگی‌ها برای بهبود عملکرد می‌پردازد. شکل (4) نشان‌دهنده نمودار جعبه‌ای این مجموعه داده می‌باشد. با توجه به شکل‌های (4) و (2) داده‌های پرت در این مجموعه داده وجود دارد.

#### 2.4. انتخاب ویژگی با الگوریتم لاسو

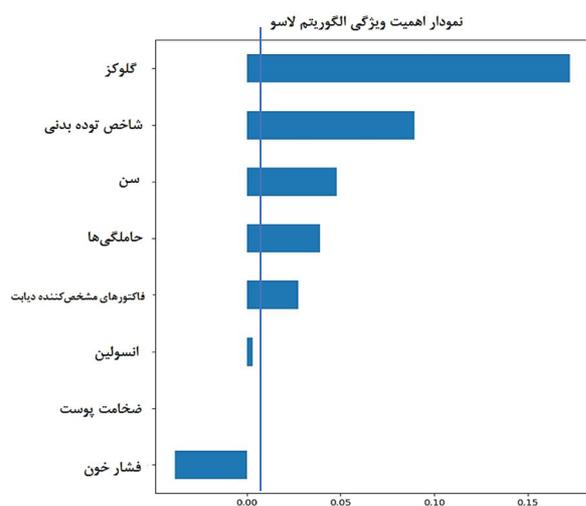
لاسو یا عملگر گزینش و انقباض کمترین قدرمطلق<sup>3</sup> در سال 1996 توسط رابرت تیب شیرانی معرفی شد [18]. این الگوریتم توسعه الگوریتم حداقل مربعات معمولی (OLS) است. روشی برای تخمین پارامترهای یک مدل رگرسیون خطی است

<sup>2</sup> Interquartile range

<sup>3</sup> Least Absolute Shrinkage and Selection Operator

<sup>1</sup> Normalization

ویژگی و هم منظم‌سازی<sup>4</sup> که برای جلوگیری از بیش‌برازش است را انجام دهد. این الگوریتم برای انتخاب ویژگی مفید است، زیرا می‌تواند مهم‌ترین ویژگی‌های در دسترس را شناسایی کند. همچنین در مساله طبقه‌بندی، لاسو اغلب برای انجام طبقه‌بندی دودویی استفاده می‌شود، جایی که هدف پیش‌بینی یکی از دو برجسب ممکن است. انتخاب ویژگی با انتخاب بهترین ویژگی‌ها و حذف ویژگی‌های نامرتبط، قدرت پیش‌بینی مدل‌ها را افزایش می‌دهد. الگوریتم لاسو، با توجه به عملکرد خوب و معمول در طبقه‌بندی دودویی و عدم بیش‌برازش، می‌تواند اهمیت ویژگی‌ها را به خوبی نشان دهد [18]. شکل (5) نتایج الگوریتم لاسو را در تشخیص اهمیت ویژگی در هشت ویژگی موجود در مجموعه داده را نشان می‌دهد. چهار ویژگی برتر که توسط الگوریتم لاسو انتخاب شده‌اند، به ترتیب شامل گلوکز، شاخص توده بدنی، سن، حاملگی‌ها، فاکتورهای مشخص‌کننده دیابت و انسولین هستند که در این پژوهش استفاده شده‌اند.



شکل (5): اهمیت ویژگی‌های موجود در مجموعه داده با استفاده از الگوریتم لاسو (منبع: محاسبات محقق).

### 3.4. طبقه‌بند

در این گام، با استفاده از ویژگی‌های انتخاب شده، به آموزش مدل پرداخته می‌شود. برای ساخت یک مدل رای‌گیری، ابتدا باید الگوریتم‌های مشارکت‌کننده در رای‌گیری آموزش داده شوند. در

که سعی می‌کند مجموع مجذور باقیمانده‌ها را به حداقل برساند. لاسو از یک تکنیک کاهش بعد استفاده می‌کند و به جای کمینه‌سازی مجموع مربعات، از یک تابع جریمه بر روی جمع قدرمطلق ضرایب مدل رگرسیونی بهره می‌برد. این تابع جریمه تعداد پارامترهای مدل را کنترل می‌کند، که این موضوع سبب کاهش پیچیدگی مدل و جلوگیری از بیش‌برازش می‌شود. انتخاب لاسو را می‌توان به صورت زیر تعریف کرد [18]:

$$\hat{\beta}^{lasso} = \arg_{\beta} \min \left\{ \frac{1}{2} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\} \quad (1)$$

در این رابطه  $\lambda$  پارامتر تنظیم‌کننده است، به این معنی که اگر مقدار آن برابر با صفر باشد، مدل به رگرسیون عادی تبدیل شده و همه متغیرها در آن حضور خواهند داشت و اگر مقدار آن افزایش یابد، تعداد متغیرهای مستقل در مدل کاهش می‌یابند. لاسو جریمه‌ای را به مجموع مجذور باقیمانده اضافه می‌کند. این جریمه برابر است با ضرب ضرایب بتا در پارامتر  $\lambda$  که در کند کردن یا شدت بخشیدن جریمه تاثیر می‌گذارد. برای مثال، اگر  $\lambda$  کمتر از 1 باشد، جریمه را کند می‌کند و اگر بالای 1 باشد، جریمه را شدت می‌بخشد. از این رو این روش انتخاب ویژگی، مجموع باقیمانده مربع‌ها را کاهش می‌دهد به شرط اینکه مجموع مقدار مطلق ضرایب کمتر از یک مقدار ثابت باشد. لاسو در ابتدا در زمینه حداقل مربعات تعریف شد، اما می‌توان آن را به مدل‌های بسیار متنوعی نیز تعمیم داد. لاسو هم دقت پیش‌بینی و هم قابلیت تفسیر مدل را بهبود می‌بخشد. اگر همبستگی بالایی در گروه پیش‌بینی‌کننده‌ها وجود داشته باشد، لاسو تنها یکی از آنها را انتخاب می‌کند و بقیه را به صفر می‌رساند. این تغییرپذیری برآوردها را با کوچک کردن برخی از ضرایب دقیقاً به صفر کاهش می‌دهد و مدل‌های قابل تفسیر تولید می‌کند.

لاسو یک مدل خطی است که با برازش ابر صفحه‌ای، داده‌ها را به دو کلاس جدا می‌کند. این مدل هم می‌تواند عمل انتخاب

<sup>4</sup> Regularization

دهد [19]. الگوریتم تقویت گرادیان شدید، با استفاده از یک تابع خطا، سعی می کند تا خطای پیش بینی را کاهش دهد. در هر مرحله، یک یادگیرنده ضعیف به مدل اضافه می شود که سعی می کند خطای باقیمانده از مراحل قبلی را کاهش دهد. این فرآیند به صورت تکرار شونده ادامه می یابد تا اینکه خطای مدل به یک حد مقبول برسد یا تعداد مراحل مشخص شده به پایان برسد. در نهایت، پیش بینی مدل نهایی، مجموع پیش بینی های تمام یادگیرندگان ضعیف خواهد بود.

### 3.3.4. کت بوست

کت بوست یک الگوریتم تجمعی یادگیری ماشین است یعنی با استفاده از مجموعه ای از یادگیرنده های ضعیف مدلی قوی تر می سازد. کاهش خطا نیز در این الگوریتم با تقویت گرادیان درخت های تصمیم در هر بار آموزش انجام می پذیرد. این الگوریتم همانند دیگر الگوریتم های تجمعی بر آن است تا با ساخت مجموعه ای از پیش بینی کننده های ضعیف و کاهش خطای آنها به یک مدل قوی دست یابد. کت بوست برخلاف دسته بندی های سنتی که نیاز به رمزگذاری یک طرفه دارند، ویژگی ها را پردازش کرده و دقت دسته بندی آن بالا است. همچنین رمزگذاری مرتب رابطه بین ویژگی ها و کلاس را در نظر می گیرد که به مراتب سبب بهبود عملکرد در مقایسه با رمزگذاری یک طرفه می شود.

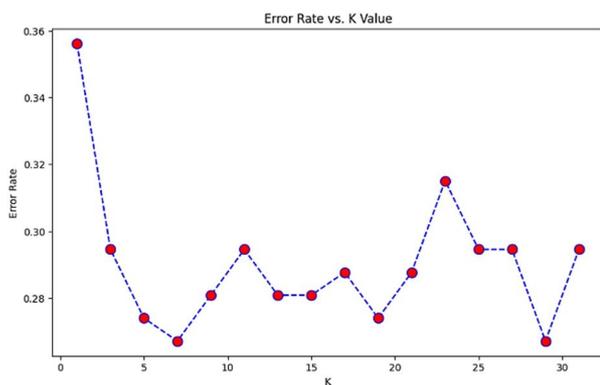
### 4.3.4. رای گیری نرم

مدل یادگیری رای گیری، مدلی است که در آن از چندین الگوریتم برای ساخت پیش بینی نهایی استفاده می شود. در این پژوهش، از الگوریتم رای گیری نرم استفاده شده است، زیرا نسبت به الگوریتم رای گیری سخت، مزایای بیشتری دارد و نیازهای این پژوهش را برآورده می کند. الگوریتم رای گیری، با توجه به مزایایی مانند جلوگیری از بیش برآزش و تفسیرپذیری پیشرفته، باعث افزایش دقت و تعمیم پذیری بیشتر الگوریتم های مورد استفاده شده است [20]. در مدل رای گیری نرم، مدل ها احتمالات برای هر کلاس ارائه می دهند و این احتمالات را برای

این پژوهش، از الگوریتم های رای گیری نرم، K-نزدیک ترین همسایه، تقویت گرادیان شدید و کت بوست استفاده شده است.

### 1.3.4. الگوریتم K-نزدیک ترین همسایه

الگوریتم K-نزدیک ترین همسایه یکی از الگوریتم های یادگیری ماشین تحت نظارت است که برای حل مسائل طبقه بندی استفاده می شود. این الگوریتم برای مسائل طبقه بندی، K-نزدیک ترین همسایه را پیدا و با اکثریت آرا نزدیک ترین همسایگان کلاس را پیش بینی می کند. شکل (6) نشان دهنده عملکرد الگوریتم K-نزدیک ترین همسایه در بازه فرد K، 1 الی 31 است. محور عمود در این شکل مقدار خطا و محور افقی مقدار K را نشان می دهد. در این پژوهش K برابر با 29 انتخاب شده است، چرا که K بزرگتر به تصمیم گیری روان تر الگوریتم منجر می شود.



شکل (6): خطای الگوریتم K-نزدیک ترین همسایه با مقادیر متفاوت K (منبع: محاسبات محقق).

### 2.3.4. الگوریتم تقویت گرادیان شدید

الگوریتم تقویت گرادیان شدید، یک تکنیک یادگیری تحت نظارت است که به خانواده الگوریتم های یادگیری ماشین مربوط به درخت های تصمیم تقویت شده با گرادیان<sup>5</sup> تعلق دارد. این الگوریتم، گروهی از یادگیرندگان ضعیف را با استفاده از فرآیند ادغام گروهی تقویت می کند و آنها را به یک یادگیرنده قدرتمند تبدیل می کند. در هر مرحله، این الگوریتم سعی می کند تا تعداد خطاهای آموزشی انجام شده توسط یادگیرنده ضعیف را کاهش

<sup>5</sup> Gradient-Boosted Decision Trees

به صورت زیر محاسبه می‌شود.

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

صحت<sup>7</sup> نشان‌دهنده آن است که پیش‌بینی مثبت مدل تا چه اندازه صحیح است و بصورت زیر محاسبه می‌شود.

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

پوشش<sup>8</sup> نشان‌دهنده این است که پیش‌بینی‌های منفی مدل تا چه اندازه صحیح است و به صورت زیر محاسبه می‌شود.

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

در روابط بالا، TP<sup>9</sup> نشان‌دهنده مثبت صحیح است، یعنی کلاس مثبت به درستی به آن نسب داده شده است، FP<sup>10</sup> نشان‌دهنده مثبت کاذب است، یعنی کلاس مثبت به اشتباه به آن نسبت داده شده است، FN<sup>11</sup> نشان‌دهنده منفی کاذب است، یعنی کلاس منفی به اشتباه به آن نسبت داده شده است و TN<sup>12</sup> نشان‌دهنده منفی صحیح است، یعنی کلاس منفی به درستی به آن نسب داده شده است.

## 2.5. نتایج

در این بخش، مجموعه داده‌های سوابق بیماران دیابت هندی پیمان به منظور پیش‌بینی دیابت مورد بررسی قرار گرفته است. با توجه به اهمیت تشخیص دیابت، در این مقاله مدلی کارآمد با استفاده از رای‌گیری نرم ایجاد شده است. روش پیشنهادی، که از رای‌گیری نرم بین سه الگوریتم K-نزدیک‌ترین همسایه، تقویت گرادیان شدید و کت‌بوست استفاده می‌کند، با ادغام نتایج این سه الگوریتم، به دقت بالاتری نسبت به هر یک از الگوریتم‌ها به تنهایی رسیده است. دقت اولیه الگوریتم‌های K-نزدیک‌ترین همسایه برابر با 81/81%، تقویت گرادیان شدید برابر با 82/46%

<sup>7</sup> Precision

<sup>8</sup> Recall

<sup>9</sup> True Positive

<sup>10</sup> False Positive

<sup>11</sup> False Negative

<sup>12</sup> True Negative

ایجاد پیش‌بینی نهایی، وزن‌دهی و میانگین‌گیری می‌کند. در مقابل، در مدل رای‌گیری سخت، نتیجه نهایی بر اساس بیشترین رای به کلاس تعیین می‌شود.

گرچه مدل یادگیری نرم ممکن است بار محاسباتی بیشتری نسبت به رای‌گیری سخت داشته باشد، اما نیاز به مدل‌ها برای تولید احتمالات یا امتیازات، اطمینان از نتیجه نهایی را فراهم می‌کند. ترکیب الگوریتم‌های یادگیری ماشین یا الگوریتم‌های رای‌گیری نرم و سخت مدلی پیشرفته برای پیش‌بینی تولید می‌کند، که پیش‌بینی‌ها را ادغام می‌کند تا نتیجه قوی‌تر و دقیق‌تری تولید شود. با استفاده از نقاط قوت مدل‌های مختلف، رای‌گیری می‌تواند نقاط ضعف آنها را کاهش دهد و بعضاً منجر به عملکرد بهتر از هر مدل تکی شود. در این پژوهش مدل رای‌گیری نرم با سه رای‌دهنده K-نزدیک‌ترین همسایه، تقویت گرادیان شدید و کت‌بوست به عنوان طبقه‌بند با هدف افزایش دقت در تشخیص دیابت مورد استفاده قرار گرفت. برای الگوریتم کت‌بوست وزنی معادل با 0/8، تقویت گرادیان شدید وزنی برابر با 0/1 و برای K-نزدیک‌ترین همسایه وزنی برابر با 0/1 انتخاب شده است. وزن‌ها به صورت آزمایشی تعیین شده و سپس بالاترین دقت حاصل از وزن 0/1-0/1-0/8 مورد استفاده قرار گرفت.

## 5. تحلیل و نتایج

در این بخش، ابتدا معیارهای ارزیابی مورد استفاده در پژوهش تعریف می‌شوند. سپس، نتایج حاصل از روش پیشنهادی بر روی مجموعه داده پیمان ارائه می‌شود. در نهایت، روش پیشنهادی با کارهای مرتبط دیگر مقایسه و تحلیل می‌شود.

### 1.5. معیارهای ارزیابی

در این بخش، معیارهای ارزیابی مورد بررسی قرار خواهد گرفت. از سه معیار ارزیابی دقت، صحت و پوشش در این پژوهش استفاده شده است. معیارهای دقت، صحت و پوشش طبق روابط زیر تعریف می‌شوند [21].

دقت<sup>6</sup> نشان‌دهنده پیش‌بینی درست خروجی توسط مدل است و

<sup>6</sup> Accuracy

داشته‌اند. از این‌رو، با توجه به مقایسه‌های انجام شده در این جدول، می‌توان گفت که مدل پیشنهادی این پژوهش، مدل کارآمدتری برای تشخیص بیماری دیابت است. با توجه به آن که کیفیت زندگی یک متغیر مهم است، از دیدگاه تحقیقات بالینی و فناوری، طراحی یک مدل یادگیری ماشین بر اساس داده‌های بیوشیمیایی که وضعیت سلامت انسان را ثبت می‌کند به یک نیاز ضروری تبدیل شده است. مدل پیشنهادی این پژوهش، که از سه الگوریتم K-نزدیک‌ترین همسایه، تقویت گرادیان شدید و کت‌بوست برای شرکت در مدل رای گیری نرم استفاده کرده است، الگوریتمی امیدبخش برای آینده تشخیص دیابت خواهد بود.

**جدول (3): مقایسه عملکرد روش پیشنهادی با دیگر پژوهش‌ها روی مجموعه داده پیمان.**

| مرجع         | نوع جداسازی داده | روش   | دقت (%) |
|--------------|------------------|---|---------|
| روش پیشنهادی | 20-80            | رای گیری نرم (با حذف داده‌های پرت)              | 94/4    |
| روش پیشنهادی | 20-80            | رای گیری نرم (بدون حذف داده‌های پرت)            | 93/6    |
| روش پیشنهادی | 20-80            | کت‌بوست   | 92      |
| [2]          | K-Fold           | ماشین بردار پشتیبان                             | 77/3    |
| [14]         | K-Fold           | LightGBM+K-NN+Adaboost                          | 90/7    |
| [15]         | K-Fold           | رای گیری نرم                                    | 85/7    |
| [16]         | K-Fold           | رای گیری سخت                                    | 77      |
| [23]         | K-Fold           | رگرسیون لجستیک                                  | 77      |
| [24]         | 30-70            | بیز اصلاح شده با الگوریتم ازدحام ذرات           | 78/6    |
| [25]         | 30-70            | جنگل تصادفی                                     | 79/57   |
| [26]         | 30-70            | درخت تصمیم                                      | 75/6    |
| [27]         | 20-80            | رای گیری نرم (جنگل تصادفی، بیز، رگرسیون لجستیک) | 79/08   |
| [28]         | K-Fold           | شبکه عصبی عمیق                                  | 89      |

و کت‌بوست برابر با 92% بوده است. نتایج روش پیشنهادی در دو حالت با حذف داده‌های پرت و بدون حذف داده‌های پرت در جدول (2) نشان داده شده است. همان‌طور که از نتایج مشخص است، با حذف داده‌های پرت، عملکرد مدل پیشنهادی به دقت 94/4%، صحت 96/5% و پوشش 92/7% رسیده است. جدول (2) نشان می‌دهد آنچه مدل پیش‌بینی کرده است تا 96% صحت دارد. همچنین مدل پیشنهادی پوشش مورد قبولی را نیز از خود نشان داده است. در این جدول نتایج حاصل از روش پیشنهادی بدون حذف مقادیر پرت نیز نشان داده شده است. این حالت، به دقت 93/6%، صحت 95/4% و پوشش 92/5% رسیده است.

**جدول (2): نتایج روش پیشنهادی برای تشخیص بیماری دیابت در 2 حالت با حذف داده‌های پرت و بدون حذف داده‌های پرت.**

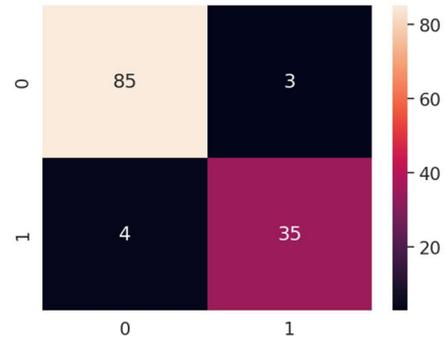
| روش          | روش رای گیری نرم | روش رای گیری نرم |
|--------------|------------------|------------------|
| حذف داده پرت | بله              | خیر              |
| جداسازی داده | 80-20            | 80-20            |
| دقت          | 94.4%            | 93.6%            |
| صحت          | 96.5%            | 95.4%            |
| پوشش         | 92.7%            | 92.5%            |

ماتریس درهم‌ریختگی مدل پیشنهادی این پژوهش با حذف داده‌های پرت در شکل (7) نشان داده شده است. ماتریس درهم‌ریختگی با نشان دادن نمونه‌هایی که به درستی و یا به غلط در هر دسته قرار گرفتند، اطلاعات مهمی در راستای کارایی مدل ارائه می‌کند [22]. همان‌طور که از ماتریس درهم‌ریختگی مشخص است روش پیشنهادی به خوبی قادر است نمونه‌های هر دو کلاس را با دقت خوبی شناسایی کند.

در جدول (3) نتایج این پژوهش با دیگر کارهای مرتبط مقایسه شده است. این جدول دقت چندین مدل مختلف برای تشخیص بیماری دیابت را نشان می‌دهد. روش پیشنهادی این پژوهش، که از روش رای گیری نرم با استفاده از سه رای‌دهنده، K-نزدیک‌ترین همسایه، تقویت گرادیان شدید و کت‌بوست استفاده می‌کند، در مقایسه با مدل‌های موجود [14]-[16]، [23]-[28] روش‌های ترکیبی، رای گیری نرم و شبکه عصبی دقت بیشتری

پژوهش تلاش شده است تا مدلی جهت افزایش دقت در تشخیص بیماری دیابت ارائه شود. برای ساخت این مدل از مجموعه داده پیما استفاده شده است. در گام نخست پیش پردازش داده‌ها انجام می‌شود، چرا که با بررسی نمونه‌های موجود در مجموعه داده می‌توان از ایجاد یک مدل دقیق اطمینان حاصل کرد. از این‌رو، در این پژوهش به حذف داده‌های پرت، جایگزین نمودن مقادیر از دست رفته با میانگین پرداخته شده است. سپس برای کاهش ابعاد و پیدا کردن ویژگی‌های مهم و تاثیرگذار، از الگوریتم لاسو استفاده شده است. در نهایت برای ساخت یک مدل با دقت بالاتر، از الگوریتم رای گیری نرم استفاده شده است. این روش، با استفاده از سه الگوریتم رای دهنده، K- نزدیک‌ترین همسایه، تقویت گرادیان شدید و کت‌بوست، دقت تشخیص بیماری دیابت را بهبود داده است.

تعارض منافع: نویسندگان اعلام می‌کنند که هیچ تعارض منافی ندارند.



شکل (7): ماتریس درهم‌ریختگی روش پیشنهادی برای تشخیص دیابت با حذف داده‌های پرت.

## 6. نتیجه گیری

دیابت یک بیماری خود ایمنی است که یکی از عوامل مهم مرگ و میر در سراسر جهان بوده و تاثیرات قابل توجهی در بیماری‌های کلیوی، بیماری‌های قلبی و از دست دادن بینایی دارد. این بیماری می‌تواند سبب امراض متعددی شود و کیفیت زندگی و سلامت انسانها را به خطر بیندازد. با بکارگیری الگوریتم‌های یادگیری ماشین در پزشکی می‌توان الگوهای پنهان بیماری، عوامل موثر در بیماری و تشخیص آن را شناسایی کرد. از این‌رو، تشخیص زودهنگام و دقیق بیماری با استفاده از این الگوریتم‌ها می‌تواند کمک شایانی به جامعه و پزشکان کند. بر این اساس، در این

## مراجع

- [1] N. Arora, A. Singh, M.Z.N. Al-Dabagh, and S.K. Maitra, "A novel architecture for diabetes patients' prediction using K-means clustering and SVM," *Math. Probl. Eng.*, vol. 2022, pp. 1-9, 2022, doi: 10.1155/2022/4815521.
- [2] D. Sisodia and D.S. Sisodia, "Prediction of diabetes using classification algorithms," *Procedia Comput. Sci.*, vol. 132, pp. 1578-1585, 2018, doi: 10.1016/j.procs.2018.05.122.
- [3] Z. Salih Ageed et al., "Comprehensive survey of big data mining approaches in cloud systems," *Qubahan Acad. J.*, vol. 1, no. 2, pp. 29-38, 2021, doi: 10.48161/qaj.v1n2a46.
- [4] W. Haoxiang and S. Smys, "Big data analysis and perturbation using data mining algorithm," *J. Soft Comput. Paradigm*, vol. 3, no. 1, pp. 19-28, 2021, doi: 10.36548/jscp.2021.1.003.
- [5] H. Wu, S. Yang, Z. Huang, J. He, and X. Wang, "Type 2 diabetes mellitus prediction model based on data mining," *Informat. Med. Unlocked*, vol. 10, pp. 100-107, 2018, doi: 10.1016/j.imu.2017.12.006.
- [6] M.M.F. Islam, R. Ferdousi, S. Rahman, and H. Bushra, "Likelihood prediction of diabetes at early stage using data mining techniques," in *Computer Vision and Machine Intelligence in Medical Image Analysis. Advances in Intelligent Systems and Computing*, vol. 992, Springer, Singapore, doi:

- 10.1007/978-981-13-8798-2\_12.
- [7] F.G. Woldemichael and S. Menaria, "Prediction of diabetes using data mining techniques," in Proc. 2nd Int. Conf. Trends Electron. Informat. (ICOEI), 2018, doi: 10.1109/icoei.2018.8553959.
- [8] C. Fiarni, E.M. Sipayung, and S. Maemunah, "Analysis and prediction of diabetes complication disease using data mining algorithm," *Procedia Comput. Sci.*, vol. 161, pp. 449-457, 2019, doi: 10.1016/j.procs.2019.11.144.
- [9] A. Aldallal and A.A.A. Al-Moosa, "Using data mining techniques to predict diabetes and heart diseases," in Proc. 4th Int. Conf. Frontiers Signal Process. (ICFSP), Poitiers, France, 2018, pp. 150-154, doi: 10.1109/ICFSP.2018.8552051.
- [10] I. Kavakiotis, O. Tsave, A. Salifoglou, N. Maglaveras, I. Vlahavas, and I. Chouvarda, "Machine learning and data mining methods in diabetes research," *Comput. Struct. Biotechnol. J.*, vol. 15, pp. 104-116, 2017, doi: 10.1016/j.csbj.2016.12.005.
- [11] A. Kumar, P. Kumar, A. Srivastava, A. Kumar, K. Vengatesan, and A. Singhal, "Comparative analysis of data mining techniques to predict heart disease for diabetic patients," in *Advances in Computing and Data Sciences (ICACDS 2020), Communications in Computer and Information Science*, vol. 1244. Springer, Singapore, 2020, doi: 10.1007/978-981-15-6634-9\_46.
- [12] T.R. Mahesh et al., "Blended ensemble learning prediction model for strengthening diagnosis and treatment of chronic diabetes disease," *Comput. Intell. Neurosci.*, vol. 2022, 2022, doi: 10.1155/2022/4451792.
- [13] A. Oza and A. Bokhare, "Diabetes prediction using logistic regression and K-nearest neighbor," in *Cong. Intell. Syst. Lect. Notes Data Eng. Commun. Technol.*, vol. 111. Springer, Singapore, 2022, doi: 10.1007/978-981-16-9113-3\_30.
- [14] M.J. Sai et al., "An ensemble of light gradient boosting machine and adaptive boosting for prediction of type-2 diabetes," *Int. J. Comput. Intell. Syst.*, vol. 16, no. 1, 2023, doi: 10.1007/s44196-023-00184-y.
- [15] A. Mahabub, "A robust voting approach for diabetes prediction using traditional machine learning techniques," *SN Appl. Sci.*, vol. 1, no. 12, 2019, doi: 10.1007/s42452-019-1759-7.
- [16] Z. Mushtaq et al., "Voting classification-based diabetes mellitus prediction using hypertuned machine-learning techniques," *Mobile Inf. Syst.*, vol. 2022, pp. 1-16, 2022, doi: 10.1155/2022/6521532.
- [17] UCI Machine Learning, "Pima Indians diabetes database," 2016, [Online]. Available: <https://www.kaggle.com/uciml/pima-indians-diabetes-database>
- [18] R. Muthukrishnan and R. Rohini, "LASSO: A feature selection technique in predictive modeling for machine learning," in Proc. IEEE Int. Conf. Adv. Comput. App. (ICACA), Coimbatore, India, 2016, pp. 18-20, doi: 10.1109/ICACA.2016.7887916.
- [19] H. Veisi, H.R. Ghaedsharaf, and M. Ebrahimi, "Improving the Performance of Machine Learning Algorithms for Heart Disease Diagnosis by Optimizing Data and Features," *Soft Comput. J.*, vol. 8, no. 1, pp. 70-85, 2019, doi: 10.22052/8.1.70 [In Persian].
- [20] F. Leon, S.-A. Floria, and C. Badica, "Evaluating the effect of voting methods on ensemble-based classification," in Proc. Int. Conf. Innovat. Intell. Syst. App. (INISTA), Gdynia, Poland, 2017, pp. 1-6, doi: 10.1109/INISTA.2017.8001122.
- [21] R. Taimourei-Yansary, M. Mirzarezaee, M. Sadeghi, and B. Nadjar Araabi, "Predicting invasive disease-free survival time in breast cancer patients using semi-supervised graph-based machine learning techniques," *Soft Comput. J.*, vol. 10, no. 1, pp. 48-69, 2021, doi: 10.22052/scj.2022.243330.1039 [In Persian].
- [22] R. Akhoondi and R. Hosseini, "A Novel Fuzzy-Genetic Differential Evolutionary Algorithm for Optimization of A Fuzzy Expert Systems Applied

- to Heart Disease Prediction,” *Soft Comput. J.*, vol. 6, no. 2, pp. 32-47, doi: 10.1001.1.23223707.1396.6.2.3.7 [In Persian].
- [23] R. Rastogi and M. Bansal, “Diabetes prediction model using data mining techniques,” *Measurement: Sensors*, vol. 24, p. 100605, 2022, doi: 10.1016/j.measen.2022.100605.
- [24] G. Battineni, G.G. Sagaro, C. Nalini, F. Amenta, and S.K. Tayebati, “Comparative machine-learning approach: A follow-up study on type 2 diabetes predictions by cross-validation methods,” *Machines*, vol. 7, no. 4, p. 74, 2019, doi: 10.3390/machines7040074.
- [25] D. Choubey, P. Kumar, S. Tripathi, and S. Kumar, “Performance evaluation of classification methods with PCA and PSO for diabetes,” *Netw. Model. Anal. Health Informat. Bioinformat.*, vol. 9, no. 1, 2020, doi: 10.1007/s13721-019-0210-8.
- [26] V. Chang, J. Bailey, Q.A. Xu, and Z. Sun, “Pima Indians diabetes mellitus classification based on machine learning (ML) algorithms,” *Neural Comput. App.*, vol. 35, pp. 16157-16173, 2023, doi: 10.1007/s00521-022-07049-z.
- [27] S. Kumari, D. Kumar, and M. Mittal, “An ensemble approach for classification and prediction of diabetes mellitus using soft voting classifier,” *Int. J. Cogn. Comput. Eng.*, vol. 2, pp. 40-46, 2021.
- [28] P. Hougue and A. G. Bigirimana, “Leveraging pima dataset to diabetes prediction: Case study of deep neural network,” *J. Comput. Commun.*, vol. 10, no. 11, pp. 15-28, 2022, doi: 10.4236/jcc.2022.1011002.