

## پیش بینی بیماری های مزمن با داده های نامتوازن توسط ماشین بردار پشتیبان گرانشی

عبدالله محمدی<sup>۱</sup>، جلال الدین نصیری<sup>۲</sup>، سهراب عفتی<sup>۳</sup>

<sup>۱</sup> دانشجوی کارشناسی ارشد علوم داده، ریاضیات کاربردی، دانشگاه فردوسی مشهد، مشهد،  
abdhmohammadi@mail.um.ac.ir

<sup>۲</sup> استادیار، گروه ریاضی کاربردی، دانشکده ریاضی دانشگاه فردوسی مشهد، مشهد،  
jnasiri@um.ac.ir

<sup>۳</sup> استاد، گروه ریاضی کاربردی، دانشکده ریاضی دانشگاه فردوسی مشهد، مشهد،  
s-effati@um.ac.ir

### چکیده

با پیشرفت تکنولوژی، روش های مبتنی بر داده برای تشخیص انواع بیماری ها، به طور گسترده ای مورد توجه قرار گرفته است. در طبقه بندی بیماری ها، تشخیص درست فرد "ناسالم" نسبت به تشخیص درست یک فرد سالم از اهمیت بیشتری برخوردار است. اغلب داده های این بیماری ها دارای جامعه ای بیمار کوچک و جامعه ای سالم بزرگتری است. در این مقاله با تعریف یک تابع وزن ویژه در مدل وزنی الگوریتم twin svm<sup>۱</sup>، نشان داده می شود اختصاص وزن به گروه کوچکتر می تواند در تشخیص طبقه ای نمونه ها موثرتر باشد. ابتدا مفاهیم پایه ای مدل را بیان نموده سپس علاوه بر روال الگوریتم های دیگر، برای داده های کلاس کوچکتر نیز وزن اختصاص داده می شود. سپس از چندین مجموعه داده ای بیماری های مزمن مانند سرطان، دیابت و آلزایمر و ... برای ارزیابی عملکرد روش استفاده نموده با مقایسه نتایج با چند روش دیگر، نشان داده می شود روش مورد استفاده می تواند با دقت بهتری نمونه ها را طبقه بندی کرده، نمونه های کلاس کوچکتر را نیز با دقت بالاتری تشخیص دهد بنابراین می توان انتظار داشت بتواند بر روش های دیگر برتری داشته باشد.

### کلمات کلیدی

ماشین بردار پشتیبان، ماشین بردار پشتیبان دوقلو، مدل وزنی، وزن گرانشی، بیماری مزمن، دیابت، آلزایمر، سرطان

### ۱- مقدمه

در مطالعه ای داده های بیماری ها یکی از موضوعات با اهمیت، طبقه بندی آن هاست در رابطه با بیماری هایی مانند سرطان و دیابت، این طبقه بندی می تواند به صورت "سالم" و "ناسالم" بیان شود. با داشتن مجموعه ای از داده های ثبت شده، یکی از روش های پیش بینی طبقه ای نمونه ای جدید استفاده از ماشین های بردار پشتیبان (SVMs) است [1]. SVM یک تکنیک قدرتمند طبقه بندی است و سعی می کند داده ها را با استفاده از دو ابرصفحه ای موازی و

ایجاد یک صفحه تصمیم گیری در بین آن ها از هم جدا کند. در بین نسخه های مختلف ارائه شده، ماشین های بردار پشتیبان دوقلو (Twin SVMs) که قید موازی بودن ابرصفحه ها را نادیده می گیرند [2] موفقیت بیشتری به دست آورده است. در کارهای نصیری و همکاران از لحاظ تئوری و هم به صورت کاربرد عملی در تشخیص حرکت انسان مورد استفاده قرار گرفته است [3], [4]. Xiong Si و Jing (2009) برای تشخیص توده در ماموگرافی دیجیتال استفاده کرده اند [5]. یکی از کاستی هایی که در svm استلندارد و twin svm وجود دارد این است که همه ای داده ها از اهمیت یکسان برخوردارند که در جهان واقعیت همیشه صادق نیست، برای غلبه بر این مشکل محققان بسیاری سعی کرده اند شیوه های مناسبی را ابداع کنند. برای طبقه بندی تومور Duan Hua و همکاران (2022) الگوریتمی با عنوان twin svm هیبریدی فازی را بکار بردند [6] که ترکیبی از مدل twin svm فازی است و با تولید یک ابرکره داده های مثبت و منفی را از هم جدا می کند. افزونه های بسیاری بر پایه ای اختصاص وزن به نمونه ها، برای این الگوریتم توسعه داده شده است که به الگوریتم های ماشین بردار پشتیبان دوقلو (weighted Twin SVMs) مشهورند [7]. در این روش ها همچنان به نسبت بین کلاس ها اهمیت داده نمی شود در حالی که ممکن است داده هایی با نسبت های نامتوازن وجود داشته باشد مانند نسبت افراد بیمار به افراد سالم. اخیراً روش هایی برای کنترل نرخ عدم تعادل کلاس<sup>۲</sup> ارائه شده است که با تعریف وزن به صورت نسبی از این نرخ، سعی می کنند این مشکل را برطرف کنند اخیراً در یک مطالعه Xiaohan Yuan و دیگران (۲۰۲۲) نیز، یک چارچوب تشخیصی جدید برای بیماری های مزمن با داده های نامتوازن ارائه کردند [8]. از دیگر روش های ارائه شده در این زمینه، الگوریتم LSFLSTSVM-CIL است که در سال ۲۰۲۲ توسط M. A. Ganaie و همکاران توسعه داده شد [9]. این شیوه سعی می کند با استفاده از نرخ عدم تعادل کلاس، یک تابع وزن برای خطای هر نمونه از داده های کلاس بزرگتر تعریف کند در این روش و روش های دیگر، میزان اهمیت خطای کلاس کوچکتر برای همه نقاط داده یکسان در نظر گرفته می شود.

<sup>2</sup> Class Imbalance Rate

<sup>1</sup> Twin Support Vector machine

## ۲-۲- ماسین بردار پشتیبان دوقلوی وزنی<sup>2</sup>

این مدل با تعریف وزن، اهمیت داده ها را در تعیین ابرصفحه ها در نظر می گیرد، با الهام از نظریه گراف، یک گراف  $K$  نزدیکترین همسایه  $G$  را برای مدل سازی ساختار هندسی محلی داده ها ایجاد می کند. ماتریس وزن  $G$  را به صورت زیر تعریف می کند:

$$W_{ij} = \begin{cases} 1 & \text{اگر } x_i \text{ در } k - \text{ همسایگی } x_j \text{ یا } x_j \text{ در } k - \text{ همسایگی } x_i \\ 0 & \text{در غیر این صورت} \end{cases} \quad (3)$$

براساس تعریف فوق برای مدل سازی فشرده درون کلاسی و تفکیک پذیری بین طبقاتی دو ماتریس گراف برای هر یک از جفت TWSVM می سازد، یک گراف درون کلاسی  $G_S$  و یک گراف بین کلاسی  $G_d$ . این دو زیرگرافهایی از  $G$  هستند.

$$W_{s,ij} = \begin{cases} 1 & \text{اگر حداقل یکی از } x_i \text{ و } x_j \text{ در همسایگی دیگری در } G_S \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases} \quad (4)$$

$$W_{d,ij} = \begin{cases} 1 & \text{اگر حداقل یکی از } x_i \text{ و } x_j \text{ در همسایگی دیگری در } G_d \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases} \quad (5)$$

ایده WLTSVM کشف اطلاعات شباهت ذاتی در نمونه های یک کلاس و استخراج بردارهای پشتیبان احتمالی موجود در نمونه های کلاس دیگر است. ماتریس وزن را با بازتعریف از روی روابط بالا به صورت زیر می سازد:

$$f_j = \begin{cases} 1 & \text{اگر وجود داشته باشد } i \text{ که } W_{d,ij} \neq 0 \\ 0 & \text{در غیر این صورت} \end{cases} \quad (6)$$

مشابه TWSVM سعی می کند دو صفحه غیر موازی بدست آورد که هر کدام با نقاط کلاس مربوطه مطابقت دارند. به دنبال تفسیر هندسی مشابه، یک مسئله بهینه سازی برای تخمین ابرصفحه کلاس ۱+ و ۱- به صورت زیر بیان می شود:

$$\text{Min } \frac{1}{2} \sum_{i=1}^{N_1} \sum_{j=1}^{N_1} W_{s,ij} (w_1^T x_j^{(1)} + b_1)^2 + C \sum_{j=1}^{N_2} \xi_j, \quad (7)$$

$$\text{s.t. } -f_j (w_1^T x_j^{(2)} + b_1) + \xi_j \geq f_j \cdot 1, \quad \xi_j \geq 0.$$

$$\text{Min } \frac{1}{2} \sum_{j=1}^{N_1} d_j (w_1^T x_j^{(1)} + b_1)^2 + C \sum_{j=1}^{N_2} \xi_j, \quad (8)$$

$$\text{s.t. } -f_j (w_1^T x_j^{(2)} + b_1) + \xi_j \geq f_j \cdot 1, \quad \xi_j \geq 0$$

که در روابط (۷) و (۸) عبارت  $f_j \cdot 1$  به معنی ضرب اسکالر  $f_j$  در برداری با درایه های ۱ است.

## ۳- روش پیشنهادی

الگوریتم LSFLTSVM-CIL<sup>3</sup> مانند مدل پایه ی Twin SVM و WLTSVM دو ابرصفحه ی غیر موازی را برای جداسازی کلاس ها از یکدیگر جستجو می کند برای کاهش خطای کلاس بزرگتر نیز در تابع هدف، از یک تابع وزن [10] برای کاهش مجموع مربعات خطای آن کلاس استفاده می کند. در اینجا به بیان خلاصه ی از این الگوریتم می پردازیم. خواننده می تواند برای جزئیات بیشتر به [9] مراجعه کند.

در این مقاله تابع وزن جدیدی بر پایه ی مفاهیم جاذبه در فیزیک بیان می شود، از این تابع در ضرایب خطای هر دو کلاس، برای به حداقل رساندن خطای تشخیص نادرست داده ها استفاده می شود، بنابراین اختصاص وزن به نمونه های کلاس کوچکتر باعث افزایش کارایی مدل می شود.

در ادامه، در بخش ۲ مفاهیم اساسی ماسین بردار پشتیبان دوقلو و وزنی معرفی می شود. در بخش ۳ روش مورد مطالعه و تابع وزن براساس مفهوم جاذبه ی گرانش بین دو جسم بیان می شود. در بخش ۴ به ارزیابی نتایج و مقایسه ی آن با روش های دیگر می پردازیم و از چندین مجموعه داده ی بیمای مزمن برای آزمایش عملکرد آن استفاده می کنیم و نتایج بدست آمده را با چند روش دیگر مقایسه می شود. این مقایسه ها نشان می دهند که این روش می تواند با دقت بالاتری بیماری های مزمن را پیش بینی کند و انتظار می رود به روش های دیگر برتری داشته باشد.

## ۲- پژوهش های گذشته

### ۲-۱- ماسین بردار پشتیبان دوقلو<sup>۱</sup>

الگوریتم Twin SVM دو ابرصفحه غیر موازی را با حل دو مسئله ی برنامه ریزی درجه ی دو به صورت زیر بدست می آورد:

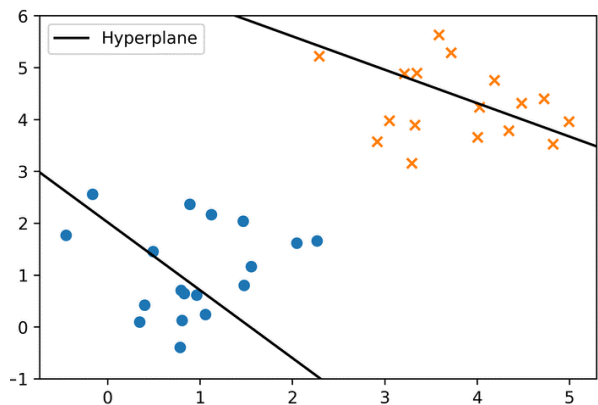
$$\text{Min } \frac{1}{2} (Aw^{(1)} + e_1 b^{(1)})^T (Aw^{(1)} + e_1 b^{(1)}) + c_1 e_2^T q \quad (1)$$

$$\text{S.t. } - (Bw^{(1)} + e_2 b^{(1)}) + q \geq e_2, \quad q \geq 0 \quad (2)$$

$$\text{Min } \frac{1}{2} (Bw^{(2)} + e_2 b^{(2)})^T (Bw^{(2)} + e_2 b^{(2)}) + c_2 e_1^T q$$

$$\text{S.t. } (Aw^{(2)} + e_1 b^{(2)}) + q \geq e_1, \quad q \geq 0$$

که در روابط (۱) و (۲)  $A$  و  $B$  به ترتیب نمونه های کلاس مثبت و منفی هستند  $b^{(i)}$  مقادیر اسکالر و بردارهای پارامترهای ابرصفحه ها،  $c_i > 0$  پارامترهای مدل،  $e_i$  بردارهایی با درایه های ۱ و  $q$  مقدار خطای ابرصفحه است. جمله ی اول در تابع هدف مجموع مربعات فاصله ی نقاط از ابرصفحه ی هر کلاس را پیاده سازی می کند و جمله ی دوم مجموع خطای متغیرها را کمینه می کند. و نمونه های جدید را براساس نزدیک بودن به ابرصفحه ها طبقه بندی می کند.



شکل ۱: ابرصفحه ی Twin SVM

<sup>3</sup> Large Scale Fuzzy Least Squares Twin SVM-Class Imbalance Learning

Twin Support Vector Machine (TSVM or Twin SVM)

<sup>2</sup> Weighted Twin Support Vector Machine (WTSVM)

### ۳-۲- حالت غیر خطی الگوریتم LSFLSTSVM-CIL

استفاده از توابع هسته<sup>۱</sup> در داده‌هایی که به صورت خطی قابل تفکیک نیستند یکی از ترفندهای مفید است [12]. برای حالت غیر خطی با انتخاب یک هسته مناسب، ماتریس  $A$  با  $K(A, A^T)$  و ماتریس  $B$  با  $K(B, B^T)$  ماتریس  $AB^T$  با  $K(A, B^T)$  و ماتریس  $BA^T$  با  $K(B, A^T)$  جایگزین می‌شود [9]. با جایگزینی موارد ذکر شده تابع تصمیم به صورت زیر اصلاح می‌شود:

$$f_1(x) = \frac{1}{c_3} \begin{bmatrix} K(x, A^t), K(x, B^t) \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} + b_1 \quad (۱۶)$$

$$b_1 = \frac{1}{c_3} (e^t \alpha + e^t \beta) \quad (۱۷)$$

$$f_2(x) = \frac{1}{c_4} \begin{bmatrix} K(x, B^t), K(x, A^t) \end{bmatrix} \begin{bmatrix} \lambda \\ \theta \end{bmatrix} + b_2$$

$$b_2 = -\frac{1}{c_4} (e^t \lambda + e^t \theta) \quad (۱۸)$$

$$\text{class}(x) = \underset{i=1,2}{\operatorname{argmin}} (|f_i(x)|)$$

و کلاس نمونه‌ی جدید مطابق رابطه‌ی (۱۸) پیش بینی می‌شود.

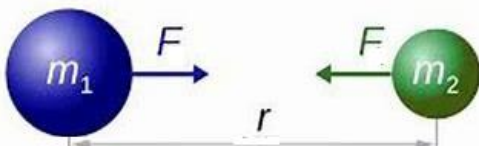
### ۳-۳- تابع وزن

در مدل‌های پیشین چون مقادیر وزن به صورت ضرایب خطای نمونه‌های کلاس بزرگتر در نظر گرفته شده‌اند فقط منجر به کاهش خطای نمونه‌های کلاس بزرگتر می‌شود و برای داده‌های کلاس کوچکتر وزنی اختصاص داده نمی‌شود، در داده‌های بیماری‌های مزمن مثل سرطان و دیابت و بیماری‌های ویروسی همه‌گیر مانند COVID-19 اهمیت بسیار زیادی وجود دارد که خطای خیلی کمتری در تشخیص کلاس کوچکتر اتفاق بیافتد در این داده‌ها تشخیص اشتباه یک نمونه به معنای این است که یک شخص ناسالم را سالم تشخیص دهیم. و این هزینه‌ی گزافی خواهد داشت. ما در اینجا ابتدا به اصلاح تابع وزن استفاده شده در مدل فوق اقدام کردیم و سپس با تخصیص وزن به داده‌های کلاس کوچکتر خطای این داده‌ها را کاهش دادیم.

### ۳-۴- تعریف تابع وزن

در فیزیک نیروی گرانش  $F$  بین دو جسم به جرم  $m_1$  و  $m_2$  با فاصله‌ی  $r$  از یکدیگر و با ثابت گرانش  $G$  به صورت زیر تعریف می‌شود:

$$F = G \frac{m_1 \times m_2}{r^2} \quad (۱۹)$$



شکل ۲- برهم کنش دو جرم در فیزیک

### ۳-۱- حالت خطی الگوریتم LSFLSTSVM-CIL

همانطور که گفته شد این الگوریتم دو ابرصفحه‌ی غیر موازی را برای جداسازی داده‌ها جستجو می‌کند که به صورت زیر داده می‌شود:

$$\begin{cases} w_1^T x + b_1 = 0 \\ w_2^T x + b_2 = 0 \end{cases} \quad (۹)$$

که در رابطه‌ی فوق  $w_1, w_2 \in \mathbb{R}^n$  و  $b_1, b_2 \in \mathbb{R}$  و  $x$  یک نمونه‌ی دلخواه از داده هاست، تابع هدف مسئله به صورت زیر تعریف می‌شود:

$$\min_{w_1, b_1, \xi_1, \eta_1} \frac{c_3}{2} (\|w_1\|^2 + b_1^2) + \frac{1}{2} \eta_1^T \eta_1 + \frac{c_1}{2} (S_2 \xi_2)^T (S_2 \xi_2) \quad (۱۰)$$

$$\text{S. t. } Aw_1 + eb_1 = \eta_1$$

$$-(Bw_1 + eb_1) + \xi_2 = e$$

$$\min_{w_2, b_2, \xi_1, \eta_2} \frac{c_4}{2} (\|w_2\|^2 + b_2^2) + \frac{1}{2} \eta_2^T \eta_2 + \frac{c_2}{2} (S_1 \xi_1)^T (S_1 \xi_1) \quad (۱۱)$$

$$\text{S. t. } Bw_2 + eb_2 = \eta_2$$

$$(Aw_2 + eb_2) + \xi_1 = e$$

که  $A$  و  $B$  به ترتیب نمونه‌های کلاس کوچکتر و بزرگتر،  $w_i, b_i$  پارامترهای ابرصفحه،  $\xi_i$  متغیرهای کمکی  $e$  بردار متشکل از درایه‌های برابر  $1, c_i$  ابرپارامترهای مدل،  $S_2$  ماتریس قطری شامل وزن خطاهای کلاس بزرگتر و  $S_1$  ماتریس همانی است.

مسئله‌های بهینه‌سازی فوق با استفاده از ضرایب لاگرانژ و اعمال شرایط K.K.T به صورت روابط (۱۲) و (۱۳) در می‌آید [9].

$$\max_{\alpha, \beta} -\frac{1}{2} (\alpha^t \ \beta^t) \hat{Q} (\alpha^t \ \beta^t)^t - c_3 \beta^t e \quad (۱۲)$$

$$\hat{Q} = \begin{bmatrix} AA^t + c_3 I & AB^t \\ BA^t & BB^t + \frac{c_3}{c_1} (S_2^{-1})^2 \end{bmatrix} + E$$

$$\max_{\lambda, \theta} -\frac{1}{2} (\lambda^t \ \theta^t) Q' (\lambda^t \ \theta^t)^t - c_4 \theta^t e \quad (۱۳)$$

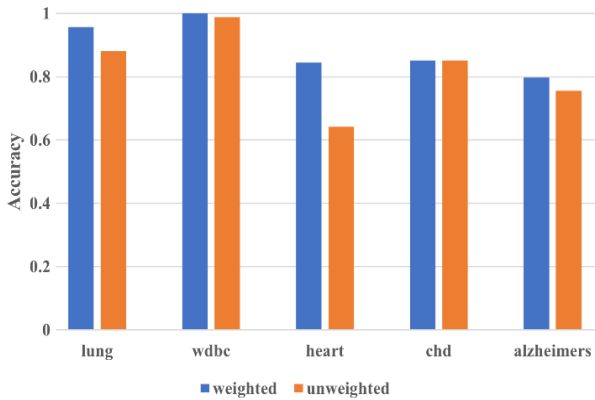
$$Q' = \begin{bmatrix} BB^t + c_4 I & BA^t \\ AB^t & AA^t + \frac{c_4}{c_2} (S_1^{-1})^2 \end{bmatrix} + E$$

و بردارهای  $w_i$  و مقادیر  $b_i$  از رابطه‌ی (۱۴) بدست می‌آید که در این رابطه  $\alpha$  و  $\beta$  مقادیر ضرایب لاگرانژ هستند.

$$\begin{bmatrix} w_1 \\ b_1 \end{bmatrix} = \frac{1}{c_3} \begin{bmatrix} A^t & B^t \\ e^t & e^t \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} \quad (۱۴)$$

پس از حل مسئله‌ی فوق تابع تصمیم‌گیری برای پیش‌بینی طبقه‌ی نمونه‌ی جدید از رابطه‌ی زیر استفاده می‌شود:

$$\text{class}(x) = \underset{i=1,2}{\operatorname{argmin}} (|x^T w_i + b_i|) \quad (۱۵)$$



شکل ۳: تفاوت دقت در دو حالت وزن دهی

شکل ۳ دیده می‌شود مدل مورد مطالعه در ۴ مورد از ۵ مورد نتایج بهتر و در یک مورد دارای نتیجه‌ی برابر بدست آورده است. می‌توان انتظار داشت وقتی به کلاس کوچکتر وزن داده شود نتایج بهتری کسب شود.

جدول ۲: پارامترهای مدل

پارامترهای مدل	محدوده‌ی مقادیر
$c, c_0$	$0.25 \times i, \quad i: 1, \dots, 10$
$c_1, c_2, c_3, c_4, \mu$	$10^i, \quad i: -5, \dots, 5$
$r$	$\begin{cases} \frac{10^{-5}}{10^i} & i: -4, -3, -2, -1, 0 \\ 0.1 + \frac{i}{10} & i: 1, \dots, 9 \end{cases}$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (22)$$

$$Recall = \frac{TP}{TP + FN} \quad (23)$$

$$Precision = \frac{TP}{TP + FP} \quad (24)$$

$$F - score = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (25)$$

$$G - mean = \sqrt{Precision \times Recall} \quad (26)$$

جدول ۳ مقایسه‌ی بین نتایج این مدل با نتایج چند نسخه از توسعه‌های الگوریتم twin svm که به طور ویژه برای داده‌های نامتوازن طراحی شده اند را روی داده‌های سرطان و دیابت نشان می‌دهد. این داده‌ها در مخزن داده‌های keel<sup>۱</sup> و Kaggle<sup>۲</sup> در دسترس است. نتایج نشان می‌دهد در داده‌های سرطان، مدل مورد مطالعه به دقت ۱۰۰٪ دست یافته و در داده‌های دیابت رتبه‌ی دوم را دارد. در رابطه با یک بیماری همه‌گیر مانند COVID-19 و یا انواع سرطان، اگر شخص سالمی بیمار تشخیص داده شود می‌توان محتاطانه تدابیر لازم را اتخاذ کرد، اما اگر شخص بیماری سالم تشخیص داده شود ممکن است عواقب جبران ناپذیری داشته باشد در این

با بهره‌گیری از این مفهوم و اندکی تغییر در آن، برای نمونه‌ی  $x$  وزن خطای آن را به صورت زیر تعریف می‌کنیم:

$$weight(x) = G \frac{N_1 \times N_2}{r(x) + 1} \quad (20)$$

که در رابطه‌ی فوق  $N_1$  تعداد نمونه‌ها در همسایگی نمونه‌ی  $x$  به همراه خود آن، و  $N_2$  نیز همین مفهوم برای مرکز کلاس به شعاع  $\epsilon$  است، در رابطه‌ی فوق  $r(x)$  فاصله‌ی نمونه‌ی  $x$  از مرکز کلاس است. این رابطه برای محاسبه‌ی عناصر قطری ماتریس‌های  $S_1$  و  $S_2$  استفاده می‌شود، این درحالی است که در مرجع [9]، [10] و [12] فقط به ماتریس  $S_2$  وزن داده می‌شود و  $S_1$  یک ماتریس همانی در نظر گرفته شده است. مقدار  $IR$  برای نقاط کلاس کوچکتر برابر یک و برای نقاط کلاس بزرگتر به صورت رابطه‌ی زیر تعریف می‌شود:

$$IR = \frac{\text{number of samples in class B}}{\text{number of samples in class A}} \quad (21)$$

در اینجا نمونه‌ی  $x$  مانند یک جسم با جرم  $N_1$  در نظر گرفته می‌شود که عضویتش در کلاس تحت تاثیر مرکز کلاس با جرم  $N_2$  می‌باشد. و از این رو متناسب با وزنش در تعیین ابرصفحه نقش ایفا می‌کند. در رابطه‌ی (۲۰) اهمیت نمونه‌ی  $x$  با توجه به تعداد نقاط اطراف و همچنین فاصله‌ی آن از مرکز کلاس سنجیده می‌شود. هر چه تعداد نقاط اطراف یک نمونه بیشتر و به مرکز کلاس نزدیک‌تر باشد دارای وزن بیشتری در این کلاس خواهد بود. می‌دانیم فاصله‌ی یک داده‌ی پرت از مرکز کلاس زیاد است علاوه بر این به دلیل دورافتادگی، در اطراف خود نقاط بسیار اندکی دارد و یا یک نقطه‌ی کاملاً تنهاست. با این توصیفات چنین نقطه‌ای وزن بسیار اندکی خواهد داشت و بنابراین از اهمیت کمتری برخوردار است.

جدول ۱: نسبت عدم تعادل کلاس

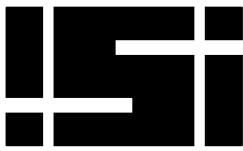
Dataset	Lung	wdbc	Heart	chd	Alzheimer
Imbalance ratio	7.0	1.6	1.23	5.52	2.33

#### ۴- نتایج داده‌های آزمایشی

در این بخش عملکرد روش خود را روی چند نمونه از داده‌های بیمه‌ای‌های مزم مورد ارزیابی قرار می‌دهیم، برای حالت غیرخطی از کرنل گوسی  $K(x_1, x_2) = e^{-\mu \|x_1 - x_2\|^2}$  و برای حالت خطی از کرنل خطی  $K(x_1, x_2) = x_1 \cdot x_2^T$  استفاده کردیم. مقادیر پارامترها از جدول ۱ و با استفاده از جستجوی تصادفی بدست آمده است، برای ارزیابی مدل از روابط (۲۲) تا (۲۶) استفاده شده است. در شکل ۳ عملکرد مدل در دو حالت مورد بررسی قرار گرفته است. یک بار به داده‌های هر دو کلاس وزن تخصیص داده شده و در حالت دوم وزن داده‌های کلاس کوچکتر برابر یک در نظر گرفته شده است. نسبت عدم تعادل این داده‌ها در جدول ۲ نشان داده شده است. در

<sup>۲</sup><https://kaggle.com>

<sup>۱</sup><https://sci2s.ugr.es/keel/development.php>



جدول ۳: مقایسه‌ی دقت مدل مورد مطالعه در نمونه داده‌های Pima و Breast Cancer با ۴ مدل دیگر

Datasets	HFSWLSTSVMS*	KWRUTSVMS-CIL*	RFLSTSVMS-CIL* ( $c_0, c_1, \mu$ )	IFW-LSTSVMS-CIL* ( $c_1, c_2, k, \mu$ )	This study ( $c, c_1, c_2, c_3, c_4, \mu, r$ )
Breast Cancer	98.55	98.01	98.76 (1.5,0.1,16)	99.07 ( $10^{-2}, 10^{-4}, 1,16$ )	<b>100</b> (0.75,2.25,10,10 <sup>2</sup> ,10 <sup>-4</sup> ,10 <sup>-5</sup> ,10 <sup>-5</sup> ,10 <sup>-5</sup> )
Pima Indian	<b>89.71</b>	71	62.77 (0.5,0.1,32)	76.77 (0.001,0.1,1,2)	79.22 (1.25,0.5,10 <sup>2</sup> ,10 <sup>2</sup> ,10 <sup>-4</sup> ,10 <sup>-3</sup> ,0.1,0.1)

\* این داده‌ها از [13]، [14] و [15] بدست آمده است.

کوچکتر در بیماریهای مزمن بسیار بیشتر از گروه بزرگتر است توجه ویژه داشتیم تا بتوانیم با افزایش دقت تشخیص صحیح در این بخش اعتبار مدل را بهبود ببخشیم، نتایج بررسی شده نشان داد که اعمال وزن به خطای گروه کوچکتر می تواند اهمیت این گروه را در تعیین طبقه بند صحیح منعکس کند با توجه به مطالب بررسی شده و نتایج موجود و مقایسه‌های انجام گرفته می توان انتظار داشت ماشین بردار پشتیبان دوقلوی با وزن گرانشی(روش مورد مطالعه) با اختصاص وزن به هر دو کلاس، در داده‌های نامتوازن عملکرد بهتری دارد. از طرفی با در نظر گرفتن مقدار نسبت عدم تعادل در گروه بزرگتر کنترل سوگیری مدل را مانند مدل‌های ذکر شده حفظ می کند به ویژه در داده‌های بیماری‌ها که اغلب اهمیت کلاس کوچکتر از کلاس بزرگتر بیشتر است این مدل عملکرد بهتری نسبت به سایر مدل‌ها دارد.

مورد اهمیت معیار recall بسیار زیاد است، این معیار با توجه به رابطه‌ی (۱۵) به معنی دقت مدل در تشخیص درست شخص واقعا بیمار است.

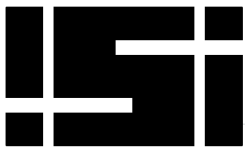
جدول ۴ عملکرد مدل را با این معیار مورد ارزیابی قرار داده است. در این جدول مشاهده می شود مدل علاوه بر نتایج مطلوب معیار صحت (Accuracy)، هماهنگ با آن نتایج معیار Recall نیز بسیار مطلوب است. در مورد Lung با میزان دقت ۹۵.۷۰٪ معیار Recall نشان دهنده‌ی این است که مدل ۹۱.۶۷٪ از بیماران را به درستی تشخیص داده است. در رابطه با داده‌های heart و hepatitis مقدار recall حتی بهتر از accuracy است و در سایر موارد نیز بیشترین تفاوت بین این دو معیار در مورد آلزایمر تقریبا ۹٪ است. این مقایسه نشان می دهد روش مورد مطالعه در هر سه مورد نتایج بهتری نشان می دهد. از نظر مقدار recall در دیتاست Pima ضعیف بوده و در hepatitis با اختلاف تقریبا ۵٪ برتری دارد.

## ۵- نتیجه گیری

در این مقاله در رابطه با داده‌های نامتوازن و اهمیت آن‌ها در الگوریتم‌های طبقه‌بندی برپایه‌ی twin svm بحث کردیم. به این نکته که اهمیت گروه

جدول ۴: عملکرد مدل و معیارهای ارزیابی

Dataset ( $c, c_1, c_2, c_3, c_4, \mu, r$ )	IR	Kernel	Accuracy	Recall	Precision	f-score	g-mean
Lung (1,1.25,1,10 <sup>3</sup> ,10 <sup>-5</sup> ,10 <sup>-3</sup> ,10 <sup>-4</sup> ,10 <sup>-5</sup> )	7	rbf	<b>0.9569</b>	<b>0.9167</b>	0.7857	0.8461	0.8487
Breast Cancer (0.75,2.25,10,10 <sup>2</sup> ,10 <sup>-4</sup> ,10 <sup>-5</sup> ,10 <sup>-5</sup> ,10 <sup>-5</sup> )	1.63	rbf	<b>1</b>	<b>1</b>	1	1	1
Pima (1.25,0.5,10 <sup>2</sup> ,10 <sup>2</sup> ,10 <sup>-4</sup> ,10 <sup>-3</sup> ,1,1)	1.98	linear	0.7705	0.7159	0.6923	0.7039	0.7040
Ckd (1,1.75,10 <sup>4</sup> ,10 <sup>5</sup> ,10 <sup>3</sup> ,10 <sup>-2</sup> ,10 <sup>-2</sup> ,0.2)	1.74	linear	<b>1</b>	<b>1</b>	1	1	1
Hepatitis (0.5,1.25,10 <sup>2</sup> ,10 <sup>2</sup> ,10 <sup>5</sup> ,10,10 <sup>-5</sup> ,10 <sup>-1</sup> ,10 <sup>-4</sup> )	3.32	linear	<b>0.9149</b>	<b>1</b>	0.6364	0.7778	0.7977
Heart (0.25,2,1,10 <sup>2</sup> ,10 <sup>2</sup> ,10 <sup>-4</sup> ,10 <sup>-3</sup> ,0.3)	1.23	linear	<b>0.8406</b>	<b>0.8689</b>	0.7910	0.8290	0.7516
Alzheimer (1.75,2,10 <sup>-2</sup> ,10 <sup>-1</sup> ,1,0.1,0.1,0.1)	2.33	rbf	0.7972	0.7025	0.6967	0.6996	0.6996



- [14] M.A. Ganaie, M. Tanveer, “KNN weighted reduced universum twin SVM for class imbalance learning”,
- [15] M. Tanveer, Senior Member, IEEE, M. A. Ganaie, A. Bhattacharjee, and C. T. Lin, “Intuitionistic Fuzzy Weighted Least Squares Twin SVMs”, IEEE.

## مراجع

- [1] V. Vapnik “Support Vector Networks”, Nature of Statistical Learning Theory (ser. Statistics for Engineering and Information Science). New York, NY, USA: Springer, 2000.
- [2] JA Nasiri, AM Mir, “An Enhanced KNN-based twin support vector machine with stable learning rules”, Neural computing and applications 32 (16), 12949-12969, 2020.
- [3] K Mozafari, JA Nasiri, NM Charkari, S Jalili, “Informatics and Computational Intelligence (ICI), 2011.
- [4] Jayadeva, R. Khemchandani, and S. Chandra, “Twin support vector machines for pattern classification”, IEEE Trans. Pattern Anal. Mach. Intell., vol. 29, no. 5, pp. 905–910, May 2007.
- [5] Xiong Si, Lu Jing, “Mass Detection in Digital Mammograms Using Twin Support Vector Machine-based CAD System”, WASE International Conference on Information Engineering, 2009.
- [6] DUAN Hua<sup>1</sup>, FENG Tong<sup>1</sup>, LIU Songning<sup>1</sup>, ZHANG Yulin<sup>1</sup>, and SU Jionglong, “Tumor Classification of Gene Expression Data by Fuzzy Hybrid Twin SVM”, Chinese Journal of Electronics Vol.31, No.1, Jan. 2022.
- [7] Qiaolin Yea, Chunxia Zhaoa, Shangbing Gaoa, Hao Zheng,” Weighted Twin Support Vector Machines with Local Information and its application”, Neural Networks, 31-39, 2012.
- [8] XiaohanYuan, Shuyu Chen, Chuan Sun & LuYuwen, “A novel early diagnostic framework for chronic diseases with class imbalance”, Scientific Reports, 12:8614, 2022.
- [9] M. A. Ganaie, M. Tanveer, Senior Member, IEEE, and Chin-Teng Lin, “Large-Scale Fuzzy Least Squares Twin SVMs for Class Imbalance Learning”, IEEE TRANSACTIONS ON FUZZY SYSTEMS, VOL. 30, NO. 11, NOVEMBER 2022.
- [10] B. Richhariya and M. Tanveer, “A robust fuzzy least squares twin support vector machine for class imbalance learning”, Appl. Soft Comput., vol. 71, pp. 418–432, 2018.
- [11] Mokhtar S. Bazaraa John J. Jarvis Hanif D “Linear Programming and Network Flows”, New York: Wiley, 1977.
- [12] Avrim Blum, John Hopcroft, and Ravindran Kannan, “Foundations of Data Science”, Thursday 4th January, 2018.
- [13] Divya Tomar and Sonali Agarwal, “Hybrid Feature Selection Based Weighted Least Squares Twin Support Vector Machine Approach for Diagnosing Breast Cancer, Hepatitis, and Diabetes”, Hindawi Publishing Corporation Advances in Artificial Neural Systems Volume 2015.