

## پیش‌بینی ابتلا به بیماریهای مزمن به کمک ماشین بردار پشتیبان دوقلو با قيود نرم

حمیده فدیشه‌ای<sup>۱</sup>، جلال‌الدین نصیری<sup>۲</sup>، سهراب عفتی<sup>۳</sup>

<sup>۱</sup> گروه ریاضی کاربردی، دانشکده علوم ریاضی، دانشگاه فردوسی مشهد، مشهد

[ha.fadishhehi@mail.um.ac.ir](mailto:ha.fadishhehi@mail.um.ac.ir)

<sup>۲</sup> استادیار، گروه ریاضی کاربردی، دانشکده علوم ریاضی، دانشگاه فردوسی مشهد، مشهد

[Jnasiri@um.ac.ir](mailto:Jnasiri@um.ac.ir)

<sup>۳</sup> استاد، گروه ریاضی کاربردی، دانشکده علوم ریاضی، دانشگاه فردوسی مشهد، مشهد

[s-effati@um.ac.ir](mailto:s-effati@um.ac.ir)

### چکیده

بیماری‌های مزمن از چالش‌های جدی مربوط به سلامت انسان می‌باشند. تشخیص بهنگام آن‌ها می‌تواند برای داشتن سبک زندگی سالم مفید باشد. روش‌های یادگیری ماشین از جمله ماشین بردار پشتیبان برای پیش‌بینی این بیماری‌ها قابل استفاده هستند. در این بررسی به کمک یک روش ماشین بردار پشتیبان دوقلو با قيود نرم، سعی شده‌است تا به کمک اطلاعات پزشکی بیماران، پیش‌بینی شود که آیا آن‌ها به بیماری‌های مزمن، مبتلا خواهند شد یا خیر.

ماشین بردار پشتیبان عادی به دنبال یافتن دو ابرصفحه موازی با بیشترین فاصله است به طوری که داده‌های دو طبقه در دو طرف آن دو قرار گیرند. در ماشین بردار پشتیبان دوقلو، دو ابرصفحه لزوماً موازی نیستند. هر یک از دو ابرصفحه نزدیکترین فاصله را به داده‌های کلاس خود داشته و از کلاس مقابل دارای یک فاصله حداقلی می‌باشند. سرعت این روش بهتر است ولی نسبت به داده‌های نویزی عملکرد خوبی ندارد. در روش پیشنهادی با استفاده از بهینه‌سازی فازی، محدودیت‌های مسأله به صورت روابط فازی در نظر گرفته می‌شوند؛ با این کار فضای شدنی مسأله بهینه‌سازی درجه دوم گسترش یافته، به نمونه‌ها اجازه تخطی از ابرصفحات داده‌شده و تأثیر داده‌های دورافتاده در تشخیص نهایی کاهش یافته‌است. الگوریتم پیشنهادی، بر روی داده‌های بالینی پزشکی، اجرا و نتایج با روش‌های مشابه مقایسه شده‌اند روش پیشنهادی عملکرد بهتری داشته است.

کرونا، از مهمترین مسأله‌های بهداشتی در جهان هستند. در سال ۲۰۱۹، حدود ۶۳٪ از مرگ و میر جهانی ناشی از بیماری‌های مزمن بوده است. رایج‌ترین این بیماری‌ها، فشار خون، قند خون، و بیماری‌های قلبی هستند که از نتایج سبک زندگی نادرست می‌باشند [9]. عوارض ناشی از آنها باعث آسیب رسیدن به قلب، کلیه، مغز و چشم‌ها می‌شود که آثاری مخرب روی کار و زندگی افراد خواهد داشت [3]. افزون بر موارد ذکر شده این بیماران در برابر بیماری‌های عفونی (مثلاً ویروس کرونا) ایمنی کمتری دارند [11]. مطابق گزارش‌های سال ۲۰۱۹، ۴۸٪ افراد مبتلا به این ویروس بیماری مزمن نیز داشته و احتمال بروز علائم شدید در آنها بیشتر است [4]. به همه موارد فوق هزینه‌های گرانی که به دولت‌ها و خانواده‌ها تحمیل می‌شوند نیز، قابل تأمل است. بنابراین تشخیص زودهنگام بیماری کمک قابل توجهی به سلامت فرد و جامعه خواهد داشت.

پس از الگوریتم‌های شبکه عصبی، در سال ۱۹۹۵ ایده جدیدتری به عنوان ماشین بردار پشتیبان توسط وینیک و همکاران مطرح شد که در آن از روش‌های بهینه‌سازی مسایل درجه دوم، برای جداسازی داده‌های دو طبقه از یکدیگر استفاده می‌کرد. مبنای کار یافتن دو ابرصفحه موازی با بیشترین فاصله ممکن بود که داده‌های دو طبقه در دو طرف این ابرصفحات قرار گیرند [10]. در طی سال‌ها انواع جدیدتر، دقیق‌تر و کاراتری از ماشین بردار پشتیبان ارائه شد. برخی از این روش‌ها در زیر بیان شده‌اند.

### ۲- پژوهش‌های پیشین

#### ۲-۱- ماشین بردار پشتیبان

فرض کنید که  $X = \{(x_i, y_i)\}_{i=1}^m$  مجموعه نمونه‌های آموزشی با اندازه  $m$  باشد که هر  $x_i \in R^n$  یک نمونه یادگیری دارای  $n$  ویژگی در فضای ورودی است،  $y_i \in \{1, -1\}$  برچسب کلاس داده  $x_i$  و  $A_{m_1 \times n}$  مجموعه نمونه‌های با برچسب  $-1$ ،  $B_{m_2 \times n}$  مجموعه نمونه‌های با برچسب  $+1$  و  $m = m_1 + m_2$  باشند، اگر داده‌ها در فضای ورودی، به صورت خطی قابل جداسازی باشند، می‌توان ابرصفحه‌ای جداساز به فرم  $w^T x + b = 0$  یافت که در آن  $w$  یک بردار  $n$  بعدی و  $b$  کمیتی عددی است. ماشین بردار

### کلمات کلیدی

یادگیری ماشین، ماشین بردار پشتیبان دوقلو، بهینه‌سازی محدب، محدودیت‌های فازی.

### ۱- مقدمه

بیماری‌های مزمن به دلیل طولانی بودن دوران ابتلا، عوارض آزاردهنده مختلف و زمینه‌سازی برای مبتلا شدن به بیماری‌های عفونی دیگر از جمله

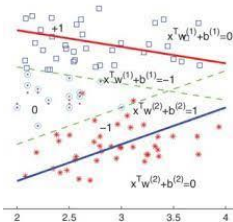
در دسته‌بندی با قیدهای نرم، محدودیت‌های مسأله به فرم محدودیت فازی درمی‌آیند:

$$y_i(\omega^T x_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, m. \quad (۷)$$

ماشین بردار پشتیبان با قیود نرم، انعطاف بیشتری نسبت به ماشین بردار پشتیبان عادی دارد؛ داده‌ها اجازه تخطی از قیدها را دارند، با این راهکار، فضای شدنی مسأله (۶) گسترش می‌یابد که در نتیجه تأثیر داده‌های نویزی کمتر می‌شود. توضیحات بیشتر در مرجع [۱] آمده‌اند.

### ۲-۳- ماشین بردار پشتیبان دوقلو

با گسترش ایده SVM در سال ۲۰۰۷، توسط جایاودا و همکاران ایده یافتن دوابرصحه غیرموازی، مطرح شد که هر یک از آن‌ها تا حد ممکن به داده‌های یک طبقه نزدیک باشند و از طبقه مقابل فاصله مطلوبی بگیرند. این الگوریتم-ها به ماشین بردار پشتیبان دوقلو<sup>۲</sup> مشهورند. در این مسایل به جای حل یک مسأله بهینه سازی درجه دوم با ابعاد بالا، دو مسأله درجه دوم جدید با ابعاد تقریباً نصف ابعاد قبل (به شرط وجود تعادل بین داده‌های دو طبقه) و با مرتبه محاسباتی  $\frac{1}{4}$  روش قبل حل می‌شوند. این بالاترین مزیت روش جدید به شمار می‌آید. گسترش‌های قابل توجهی از ماشین بردار پشتیبان دوقلو مطرح [7] و کاربردهای متعددی از این روش ارائه شده است که از جمله آن‌ها می‌توان به تشخیص رفتار انسانی [6] اشاره کرد.



شکل ۲: ماشین بردار پشتیبان دوقلو

در توضیح الگوریتم ماشین بردار پشتیبان دوقلو به فرض مجموعه‌ای از نمونه‌ها و اطلاعاتی درمورد آنها موجود است که به هر نمونه یک برچسب اختصاص داده شده و به کمک برچسب‌ها (یا ۱-۱) میتوان کل نمونه‌ها را به دو گروه دسته‌بندی کرد. بهتر است که توزیع داده‌ها متعادل باشد یعنی با نسبت تقریباً مساوی، با برچسب ۱ یا -۱ دسته‌بندی شده باشند. در روش بردار پشتیبان دوقلوبه جای یافتن دو ابرصحه موازی هم، ابرصحه‌های جداساز، ما متقاطع اند. (مثلاً  $h_1$  و  $h_2$ )

$$\begin{aligned} h_1 : x^T \omega_1 + b_1 &= 0 \\ h_2 : x^T \omega_2 + b_2 &= 0 \end{aligned} \quad (۸)$$

که در رابطه‌های فوق  $\omega_1$  و  $\omega_2$  نمایشگر مختصات ابر صفحه و  $b_1$  و  $b_2$  با یاس می‌باشند. تا حد امکان به کلاس +۱ نزدیک و از کلاس -۱ دور است و  $h_2$  تا حد ممکن به کلاس -۱ نزدیک و از کلاس +۱ دور می‌باشد. داده‌های با برچسب منفی بایدفاصله‌ای مطلوب و حداقلی (مثلاً یک واحد) از  $h_1$  داشته باشند در غیراین صورت خطای دسته‌بندی ایجاد خواهد شد و لازم است که پارامترهای مشخص‌کننده جداساز طوری تنظیم شوند که داده‌هایی از این دست کمتر رخ دهند. مشابه در مورد داده‌های دسته دوم هم چنین وضعیتی رخ خواهد داد. از این رو دو مسأله مینی‌موم سازی مطرح خواهد شد و هر یک از این دو مسأله نیاز به کمینه کردن دو مقدار متفاوت دارند؛ کم کردن فاصله ابرصحه جداساز  $h_1$  از کلاس +۱ و کم کردن تعداد نمونه‌هایی از کلاس -۱ که به  $h_1$  نزدیک شده‌اند. تمام مباحث

پشتیبان به جای یافتن یک ابرصحه جداکننده، به دنبال یافتن دوابرصحه موازی است که تا حد امکان از یکدیگر دور بوده و داده‌های دو طبقه در دو طرف آنها واقع باشند. با مدل سازی مسأله بالا و حل مسأله بهینه‌سازی (۱)، میتوان ابرصحه جداساز را مشخص کرد.

$$\min \frac{1}{2} \|\omega\|^2 \quad (۹)$$

$$\begin{aligned} S. t. \quad & y_i(\omega^T x_i + b) \geq 1, \quad i = 1, 2, \dots, m. \end{aligned}$$

نمایش ماتریسی مسأله (۱) به فرم (۲) خواهد بود:

$$\min \frac{1}{2} \|\omega\|^2 \quad (۱۰)$$

$$\begin{aligned} S. t. \quad & A\omega + e_1 b \geq e_1, \\ & B\omega + e_2 b \leq -e_2 \end{aligned}$$

در (۲)،  $e_1$  و  $e_2$  بردارهای با درایه‌های واحد و متناسب با  $A$  و  $B$  هستند. در حالتی که داده‌ها در فضای ورودی، به صورت خطی جداناپذیر هستند، مسأله SVM با حاشیه نرم به صورت (۳) مطرح خواهد شد:

$$\min \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m \xi_i \quad (۱۱)$$

$$\begin{aligned} S. t. \quad & y_i(\omega^T x_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, \quad i = 1, 2, \dots, m. \end{aligned}$$

که در آن  $\xi_i$  متغیر لغزش برای کاهش تأثیر داده نویزی  $x_i$  است. پس از نوشتن تابع لاگرانژ و در نظر گرفتن شرایط مکمل زاید و حل مسأله (۳) ابرصحه جداساز معلوم خواهد شد.

$$\min \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{k=1}^m \alpha_k \quad (۱۲)$$

$$\begin{aligned} S. t. \quad & \sum_{i=1}^m \alpha_i y_i = 0, \\ & 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, m. \end{aligned}$$

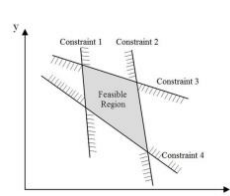
که ضرایب لاگرانژ هستند و تابع تصمیم  $D(x)$  طبق فرمول (۴) موقعیت داده جدید  $x$  را نسبت به ابرصحه مشخص خواهد کرد:

$$\begin{aligned} D(x) &= \text{sign}(\omega^T x + b) \\ &= \text{sign}(\sum_{i \in S} \alpha_i y_i x_i^T x + b) \end{aligned} \quad (۱۳)$$

در رابطه (۱۳)،  $S$  نشان‌دهنده مجموعه اندیس‌های بردارهای پشتیبان است. اگر نمونه‌ها در فضای ورودی جداناپذیر خطی نباشند، نمونه‌ها با تکنیک حقه<sup>۲</sup> در فضای ویژگی با ابعاد بیشتر نگاشته می‌شوند. در این فضا، ماهیت خود داده، اهمیت ندارد بلکه آنچه مهم است فاصله داده‌ها از یکدیگر است. با این روش ماشین بردار پشتیبان یک ابرصحه جداکننده بهینه برای جداسازی نمونه‌ها پیدا می‌کند.

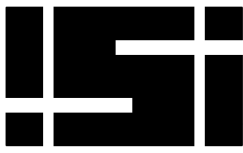
### ۲-۲- ماشین بردار پشتیبان با قیود نرم

مسأله بهینه‌سازی (۶)، برنامه‌ریزی فازی نامتقارن نامیده می‌شود، اگر محدودیت‌هایش به صورت فازی بیان شوند. برای حل این مسأله روش‌هایی پیشنهاد شده است [۲].



شکل ۳: فضای شدنی برای مسأله (۶)

$$\begin{aligned} & \text{Optimize } f(x) \\ S. t. \quad & \begin{cases} \text{constraint 1} \\ \text{constraint 2} \\ \text{constraint 3} \\ \text{constraint 4} \end{cases} \end{aligned}$$



توضیحات بیشتر در مورد داده‌هایی که به‌طور خطی جداناپذیرند، در مرجع [8] ذکر شده‌است.

گفته‌شده در مورد جداساز  $h_2$  هم صدق خواهند کرد. با توجه به اینکه مسألهٔ بهینه‌سازی درجه دوم (QPP) به فرم استاندارد زیر است:

$$\text{minimize} \quad \left(\frac{1}{2}\right)x^T P x + q^T x \quad (9)$$

$$S. t. \quad \begin{cases} Gx \leq h \\ Ax = b \end{cases}$$

دو مسألهٔ بهینه‌سازی به صورت زیر تعریف می‌شوند:

### ۳- راه حل پیشنهادی: ماشین بردار پشتیبان دوقلو با قيود نرم<sup>۵</sup>

#### ۳-۱- طبقه‌بند خطی

در این بررسی تلاش خواهد شد تا به کمک یک نوع ماشین بردار پشتیبان دوقلو با محدودیت‌های نرم، ابتلائی فرد به بعضی از انواع بیماری‌های مزمن، پیش‌بینی شود. روش پیشنهادی، در مقایسه با ماشین بردار پشتیبان دوقلو، هنگام مواجهه با داده‌های نویزی عملکرد بهتری دارد.

در روش پیشنهادی سعی شده‌است تا با ایجاد تغییر در محدودیت‌های مسألهٔ بهینه‌سازی درجه دوم در ماشین بردار پشتیبان دوقلو، به‌جای نامساوی‌های ساده از نامساوی‌های فازی استفاده شود. مسأله‌های بهینه‌سازی جدید از قرار زیر هستند:

$$\text{min} \quad \frac{1}{2} \|A\omega_1 + e_1 b_1\|^2 + C_1 e_2^T q_1 \quad (17)$$

$$S. t. \quad -(B\omega_1 + e_2 b_1) + q_1 \geq e_2, q_1 \geq 0$$

$$\text{min} \quad \frac{1}{2} \|B\omega_2 + e_2 b_2\|^2 + C_2 e_1^T q_2 \quad (18)$$

$$S. t. \quad (A\omega_2 + e_1 b_2) + q_2 \geq e_1, q_2 \geq 0$$

علامت  $\geq$  به این معنی است که برای نمونه‌های یادگیری اجازه‌ی تخطی از حدود هم داده می‌شود. در واقع با این کار فضای شدنی مسأله گسترش می‌یابد و جواب‌های بهینه از فضای بزرگتری می‌توانند انتخاب شوند.

ابتدا می‌بایست دو مسأله‌ی بهینه‌سازی (۱۷) و (۱۸) که مسأله‌های بهینه‌سازی غیرخطی درجهٔ دوم با تابع هدف محدب و قيود نامساوی فازی هستند، حل شوند. در گام اول برای حل مسأله، محدودیتها را از فرم ماتریسی خارج کرده و برای هر نمونه از هر طبقه  $i=1$ ، یک محدودیت به شکل (۱۹) در نظر گرفته می‌شود.

$$-(\omega_1^T x_i + b_1) \geq 1 - q_{1i}, \quad q_{1i} \geq 0 \quad (19)$$

$$i = 1, \dots, m_2.$$

که  $x_i$ ،  $i$  - امین نمونه از داده‌های با برچسب  $-1$  و  $q_{1i}$  خطای مربوط به این نمونه در اثر تخطی از محدودیت نظیر و نزدیک شدن به  $h_1$  می‌باشد. به‌طور مشابه برای نمونه  $i$  - ام از طبقهٔ  $+1$  محدودیت به شکل (۲۰) خواهد بود:

$$(\omega_2^T x_i + b_2) \geq 1 - q_{2i}, \quad q_{2i} \geq 0 \quad (20)$$

$$i = 1, \dots, m_1.$$

که  $x_i$ ،  $i$  - امین نمونه از داده‌های با برچسب  $+1$  و  $q_{2i}$  خطای مربوط به این نمونه در اثر تخطی از محدودیت نظیر و نزدیک شدن به  $h_2$  می‌باشد. برای نمونه‌های با برچسب  $+1$  طرفین محدودیت نظیر در  $(+1)$  و برای نمونه‌های با برچسب  $-1$  طرفین محدودیت نظیر در  $(-1)$  ضرب شده و در نهایت محدودیتها به شکل (۲۱) ادغام خواهند شد.

$$y_i(\omega_1^T x_i + b_1) \geq 1 - q_{1i}, \quad (21)$$

$$q_{1i} \geq 0, \quad i = 1, \dots, m_2.$$

$$\text{min} \quad \frac{1}{2} \|A\omega_1 + e_1 b_1\|^2 + C_1 e_2^T q_1 \quad (10)$$

$$S. t. \quad -(B\omega_1 + e_2 b_1) + q_1 \geq e_2, q_1 \geq 0$$

$$\text{min} \quad \frac{1}{2} \|B\omega_2 + e_2 b_2\|^2 + C_2 e_1^T q_2 \quad (11)$$

$$S. t. \quad (A\omega_2 + e_1 b_2) + q_2 \geq e_1, q_2 \geq 0$$

که در آن،  $q_2 \in \mathcal{R}^{m_1}$ ،  $q_1 \in \mathcal{R}^{m_2}$ ،  $C_1$  و  $C_2$  پارامترهای خطا با مقادیر مثبت و سایر متغیرها قبلاً تعریف شده‌اند. برای حل مسأله‌های فوق، بعد از نوشتن دوگان لاگرانژ و شرایط مکمل زاید<sup>۶</sup>، به دو مسألهٔ جدید زیر می‌رسیم:

$$\text{min} \quad \frac{1}{2} \alpha^T G (H^T H + \epsilon I)^{-1} G^T \alpha - e_2^T \alpha \quad (12)$$

$$S. t. \quad 0 \leq \alpha \leq C_1 e_2$$

$$\text{min} \quad \frac{1}{2} \beta^T H (G^T G + \epsilon I)^{-1} H^T \beta - e_1^T \beta \quad (13)$$

$$S. t. \quad 0 \leq \beta \leq C_2 e_1$$

پس از بهینه‌سازی و یافتن بردارهای  $\alpha$  و  $\beta$ ، ابر صفحه‌های دو کلاس از دستورات زیر به دست می‌آیند:

$$\begin{bmatrix} \omega_1 \\ b_1 \end{bmatrix} = -(H^T H + \epsilon I)^{-1} G^T \alpha \quad (14)$$

$$\begin{bmatrix} \omega_2 \\ b_2 \end{bmatrix} = (G^T G + \epsilon I)^{-1} H^T \beta \quad (15)$$

در رابطه (۱۵) ماتریسهای  $H$  و  $G$  به صورت  $H = [A \ e_1]$  و  $G = [B \ e_2]$  تعریف شده‌اند. مقدار کوچک و مثبت  $\epsilon$  به قطر اصلی دوماتریس افزوده شده تا ماتریسها وارون‌پذیر باشند. در نهایت با حل مسأله فوق و یافتن ابرصفحات جداکننده، فاصلهٔ نمونه‌های جدید از هر ابر صفحه محاسبه شده و مینی‌موم فاصله، مشخص خواهد کرد که نمونه ذکر شده به کدام دسته تعلق دارد.

در صورتیکه داده‌ها بصورت خطی قابل تفکیک نباشند، شبیه SVM از حقه‌ی هسته استفاده شده و ابر صفحه‌هایی با معادله‌های زیر تشکیل خواهند شد:

$$h_1 : K(x^T, C^t)\omega_1 + b_1 = 0 \quad (16)$$

$$h_2 : K(x^T, C^t)\omega_2 + b_2 = 0$$

های مسأله یعنی  $(\omega_1, b_1, q_1) \in R^{m_2+1+n}$  از این فضا انتخاب می‌شوند. هنگام حل مسأله این شرایط باید برآورده شوند. استفاده از نامساوی فازی به

این محدودیت‌ها انعطاف بیشتری می‌بخشد. پارامترهای  $\alpha$  و  $d_{1i}$  توسط کاربر تعیین می‌شوند، پارامتر  $\alpha$ ، سطحی را مشخص می‌کند که برای سطوح پایین‌تر از آن میزان تعلق فازی برای نامساوی موجود در قیود (۲۱)، صفر است. هر چه این مقدار به عدد ۱ نزدیک‌تر باشد، قیود مسأله به TWSVM نزدیک‌تر می‌شود.  $d_{1i}$  ها (که در اینجا مقدارشان مساوی  $d_1$  است)، مقدار تحملی هستند که می‌توان به نمونه‌ها نسبت داد. هر چه مقدار تحملی که نسبت می‌دهیم بیشتر باشد، نمونه‌های کلاس ۱- می‌توانند به ابرصفحه جداساز کلاس ۱+ نزدیک‌تر شوند. این بدان معنی است که تعداد بردارهای پشتیبان بیشتر شده و داده‌های بیشتری در تعیین ابرصفحه بهینه نقش خواهند داشت. تعریف می‌کنیم:

$$\mu_i: R^{m_2+1+n} \rightarrow (0,1], \quad i = 1, 2, \dots, m_2.$$

$$\mu_i(\omega_1, b_1, q_1) = \begin{cases} 1, & \text{if } y_i(\omega_1^T x_i + b_1) \geq 1 - q_{1i} \\ \frac{y_i(\omega_1^T x_i + b_1) - 1 + q_{1i} + d_{1i}}{d_{1i}}, & \text{if } 1 - (q_{1i} + d_{1i}) \leq y_i(\omega_1^T x_i + b_1) \leq 1 - q_{1i} \\ 0, & \text{if } y_i(\omega_1^T x_i + b_1) \leq 1 - (q_{1i} + d_{1i}) \end{cases}$$

زیرا:

$$\frac{y_i(\omega_1^T x_i + b_1) - 1 + q_{1i} + d_{1i}}{d_{1i}} \geq \alpha \Rightarrow \quad (30)$$

$$y_i(\omega_1^T x_i + b_1) \geq 1 - q_{1i} - d_{1i}(1 - \alpha)$$

در نهایت مدل RTWSVM عبارت است از:

$$\text{Minimize } \frac{1}{2} \|A\omega_1 + e_1 b_1\|^2 + C_1 e_2^T q_1 \quad (31)$$

$$\text{S.t. } -(B\omega_1 + e_2 b_1) \geq e_2 - q_1 - d_1(1 - \alpha)$$

$$q_{1i} \geq 0, \quad i = 1, \dots, m_2.$$

گام بعدی حل مسأله بهینه‌سازی غیرخطی درجه دوم با قیدهای غیرفازی  $\gamma$  است. تابع هدف اولیه همراه با محدودیت‌هایش، به تابع لاگرانژ تبدیل می‌شوند:

$$Q(\omega_1, b_1, q_1, \beta_1, \gamma_1) = \frac{1}{2} \|A\omega_1 + e_1 b_1\|^2 + \quad (32)$$

$$C_1 e_2^T q_1 - \beta_1^T \{-(B\omega_1 + e_2 b_1) + q_1 + d_1(1 - \alpha) - e_2\} - \gamma_1^T q_1$$

که در آن  $\gamma_1 = (\gamma_{11}, \dots, \gamma_{1m_2})^T$ ,  $\beta_1 = (\beta_{11}, \dots, \beta_{1m_2})^T$  ضرایب لاگرانژ و نامنفی می‌باشند.

$$\frac{\partial Q(\omega_1, b_1, q_1, \beta_1, \gamma_1)}{\partial \omega_1} = 0 \Rightarrow \quad (33)$$

$$A^T(A\omega_1 + e_1 b_1) + B^T \beta_1 = 0$$

$$\frac{\partial Q(\omega_1, b_1, q_1, \beta_1, \gamma_1)}{\partial b_1} = 0 \Rightarrow \quad (34)$$

$$e_1^T (A\omega_1 + e_1 b_1) + e_2^T \beta_1 = 0$$

$$\frac{\partial Q(\omega_1, b_1, q_1, \beta_1, \gamma_1)}{\partial q_1} = 0 \Rightarrow \quad (35)$$

$$C_1 e_2^T - \beta_1^T - \gamma_1^T = 0$$

به‌علاوه باید شرایط مکمل زاید نیز برقرار شوند:

$$\beta_1^T \{-(B\omega_1 + e_2 b_1) + q_1 + d_1(1 - \alpha) - e_2\} = 0 \quad (36)$$

$$e_2^T \beta_1 = 0, \quad \gamma_1^T q_1 = 0$$

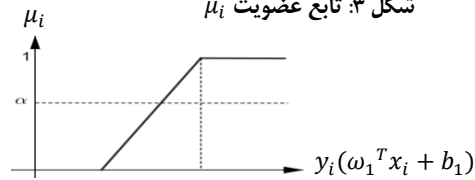
از ترکیب روابط (۳۳) و (۳۴):

$$y_i(\omega_2^T x_i + b_2) \geq 1 - q_{2i},$$

$$q_{2i} \geq 0, \quad i = 1, \dots, m_1.$$

گام دوم تعیین توابع عضویت خطی  $\mu_i, \mu'_i$  برای نامساوی‌های فازی موجود در قیود مسأله است. با توجه به شکل (۲)، در مورد اولین محدودیت در رابطه (۲۱) و با تمرکز بر داده‌های کلاس ۱-، تابع عضویت  $\mu_i$  مربوط به این طبقه تفصیل بررسی خواهد شد.  $\mu'_i$  مشابه خواهد بود.

شکل ۳: تابع عضویت  $\mu_i$



$$1 - (q_{1i} + d_{1i}) \quad 1 - q_{1i}$$

طبق رابطه (۲۱) هریک از نمونه‌های فوق، یک محدودیت به مسأله

تحلیل می‌کنند. مجموعه این محدودیت‌ها فضای شدنی را می‌سازد. متغیر-

(۲۲)

و برای هر قید (۲۱) تعریف می‌کنیم:

$$X_i = \{(\omega_1, b_1, q_1) \in R^{m_2+1+n} \mid y_i(\omega_1^T x_i + b_1) \geq 1 - q_{1i}, q_{1i} \geq 0, \quad i = 1, \dots, m_2\} \quad (23)$$

آن‌گاه جواب‌های مسأله باید از مجموعه  $X$  با تعریف زیر انتخاب شوند:

$$X = \bigcap_{i \in I} X_i, \quad I = \{1, 2, \dots, m_2\} \quad (24)$$

بنابراین مسأله  $rtwsvmI$  را می‌توان به صورت زیر نوشت:

$$\text{Minimize } \{Q(\omega_1, b_1, q_1) = \quad (25)$$

$$\frac{1}{2} \|A\omega_1 + e_1 b_1\|^2 + C_1 e_2^T q_1 \mid (\omega_1, b_1, q_1) \in X\}$$

برای حل (۲۵) از روش برش آلفا استفاده می‌شود. این مسأله در شکل ۳ دیده می‌شود. جواب‌های (۲۵) باید از مجموعه

$$X_\alpha = \{(\omega_1, b_1, q_1) \in R^{m_2+1+n} \mid \mu_X(\omega_1, b_1, q_1) \geq \alpha\} \quad (26)$$

انتخاب شوند که

$$\mu_X(x) = \inf \{\mu_i(x), i \in I\} \quad (27)$$

مجموعه  $X_\alpha$  آلفا برش قید  $i$  ام است و  $\alpha \in (0, 1]$ .

جواب بهینه مسأله (۱۷) با  $\alpha$  داده شده برابر است با:

$$S(\alpha) = \quad (28)$$

$$\begin{aligned} & \{(\omega_1, b_1, q_1) \in R^{m_2+1+n} \mid \\ & \frac{1}{2} \|A\omega_1 + e_1 b_1\|^2 + C_1 e_2^T q_1 = \\ & \text{Min } \frac{1}{2} \|A\omega'_1 + e_1 b'_1\|^2 + \\ & C_1 e_2^T q'_1, (\omega'_1, b'_1, q'_1) \in X_\alpha\} \end{aligned}$$

پس می‌توان نوشت:

$$X_\alpha = \bigcap_{i \in I} \{(\omega_1, b_1, q_1) \in R^{m_2+1+n} \mid \quad (29)$$

$$y_i(\omega_1^T x_i + b_1) \geq r_{1i}(\alpha), q_{1i} \geq 0,$$

$$i = 1, \dots, m_2\}$$

که  $r_{1i}(\alpha) = 1 - q_{1i} - d_{1i}(1 - \alpha)$

$$U_1 = [\omega_1 b_1]^T \text{ و}$$

مشابه نسخه خطی کلاس نمونه جدید  $x$  با محاسبه فاصله عمودی آن از دو ابرسطح مشخص می شود؛ به طوری که تابع تصمیم برای نسخه غیرخطی به صورت زیر تعریف می گردد:

$$\underset{j=1,2}{\operatorname{argmin}} |K(x^T, C^T)\omega_j + b_j| \quad (45)$$

#### ۴- نتایج و ارزیابی

##### ۴-۱- پایگاه داده

داده های مورد استفاده در این تحقیق، از پایگاه های داده UCI و Kaggle استخراج شده اند. به دلیل اهمیت موضوع چهار مجموعه داده: CKD، PIDDD، Hep و WDBC که به ترتیب مرتبط با بیماری کلیه، هپاتیت، دیابت و سرطان سینه هستند، مورد بررسی قرار گرفته اند. داده های توصیفی از دست رفته با اولین مد در ستون مزبور و داده های عددی با صفر جایگزین شده اند. جدول (۱) مجموعه داده را توصیف می کند.

مجموعه داده	تعداد نمونه ها	تعداد نمونه های مثبت	تعداد نمونه های منفی	تعداد ویژگی از دست رفته	داده
CKD	۴۰۰	۳۴۸	۱۵۲	۲۴	ندارد
Hep	۱۵۵	۳۲	۱۲۳	۱۹	دارد
PIDDD	۷۶۸	۲۶۸	۵۰۰	۸	ندارد
WDBC	۵۶۹	۲۱۲	۳۵۷	۳۰	دارد

جدول ۱: مشخصات مجموعه داده برای ارزیابی روش RTWSVM

##### ۴-۲- نحوه پیاده سازی و اجرای الگوریتم ها

تمام روش ها در زبان برنامه نویسی پایتون در محیط آنلاین google colab پیاده سازی شده است. آزمایش ها روی یک کامپیوتر شخصی با پردازنده core i5، سیستم عامل ویندوز ۱۰ و ۸ گیگابایت حافظه صورت گرفته است. از کتابخانه cvxopt برای بهینه سازی استفاده شده است. از آن جا که در دنیای واقعی داده ها جداپذیر خطی نیستند، از تابع هسته گوسی که قدرت تعمیم پذیری بهتری دارد استفاده شده است:

$$K(x, x') = \exp\left(-\frac{\|x-x'\|^2}{2\gamma^2}\right) \quad (46)$$

دقت روش RTWSVM به انتخاب پارامترهای بهینه بسیار وابسته است. بدین منظور از روش جستجوی شبکه ای برای پیدا کردن پارامتر بهینه استفاده شده است. پارامترهای خطا از مجموعه  $\{2^i | i = -7, -6, \dots, 7\}$  و پارامترهای  $\alpha_1, \alpha_2, d_1, d_2$  که به ترتیب، برای سطح برش آلفا و ضریب تحمل الگوریتم، مورد استفاده بوده اند از مجموعه  $[0.۰۹۵, 0.۰۷۵, 0.۰۲۵]$  و پارامتر تابع هسته  $\gamma$  نیز از مجموعه  $\{2^i | i = -15, -6, \dots, 5\}$  انتخاب شده اند. برای تعیین میزان دقت روش، از اعتبارسنجی متقابل استفاده شده است؛ الگوریتم ارائه شده با ۹۰٪ داده ها ( $k=10$ ) آموزش دیده و با ۱۰٪ باقیمانده آزمایش شده، این مراحل ده بار تکرار شده و پس از ۱۰ بار تکرار، میانگین و انحراف معیار دقت روش پیشنهادی به دست آمده اند. شکل ۵ بیانگر تأثیر تغییرات پارامتر گاما بر میزان دقت روش پیشنهادی در پایگاه داده CKD است.

$$U_1 = -(H^T H + \epsilon I)^{-1} G^T \beta_1 \quad (37)$$

که در آن:  $U_1 = [\omega_1 b_1]^T$  و  $G = [B e_2]$ ،  $H = [A e_1]$  جایگزاری (۳۷) در (۳۲) و ساده کردن روابط و استفاده از رابطه (۳۵)، مسأله نهایی به فرم زیر خواهد بود:

$$\underset{\beta_1}{\operatorname{Min}} \frac{1}{2} \beta_1^T G (H^T H + \epsilon I)^{-1} G^T \beta_1 - \beta_1^T \{d_1(1-\alpha) - e_2\} \quad (38)$$

$$S. t. \quad 0 \leq \beta_1 \leq C_1 e_2$$

مشابه تمام روابط (۲۲) تا (۳۶) برای داده های کلاس دیگر نیز برقرار است. با حل (۳۸) و به دست آمدن ضرایب لاگرانژ و مشخص شدن دو ابرصفحه غیرموازی، تابع تصمیم (۳۹) مشخص خواهد کرد که نمونه جدید به کدام طبقه تعلق می گیرد.

$$\underset{j=1,2}{\operatorname{argmin}} |x^T \omega_j + b_j| \quad (39)$$

##### ۳-۲- طبقه بند غیرخطی

توسعه روش پیشنهادی RTWSVM برای حالت غیرخطی مشابه SVM استاندارد است. برای این حالت، داده های آموزش با استفاده از تابع  $\phi$  به فضای ویژگی نگاشت می شوند. با تعریف تابع هسته:

$$K(x_i, x_j) = \phi^T(x_i)\phi(x_j) = z_i \cdot z_j$$

ابرسطح های جداکننده به فرم (۱۶) تبدیل می شوند. در فضای ویژگی مسأله اولیه به صورت (۴۰) و (۴۱)

$$\underset{(\omega_1, b_1)}{\operatorname{min}} \frac{1}{2} \|K(A, C^T)\omega_1 + e_1 b_1\|^2 + C_1 e_2^T q_1 \quad (40)$$

$$S. t. \quad - (K(B, C^T)\omega_1 + e_2 b_1) + q_1 \geq e_2, \\ q_1 \geq 0$$

$$\underset{(\omega_2, b_2)}{\operatorname{min}} \frac{1}{2} \|K(B, C^T)\omega_2 + e_2 b_2\|^2 + C_2 e_1^T q_2 \quad (41)$$

$$S. t. \quad (K(A, C^T)\omega_2 + e_1 b_2) + q_2 \geq e_1, \\ q_2 \geq 0$$

که در آن  $C = [A B]$  و  $K$  تابع هسته دلخواه می باشد.

مسأله برای (۴۰) ادامه خواهد یافت. برای (۴۱) روش ها مشابه خواهد بود.

$$\underset{(\omega_1, b_1)}{\operatorname{min}} \frac{1}{2} \|K(A, C^T)\omega_1 + e_1 b_1\|^2 + C_1 e_2^T q_1 \quad (42)$$

$$S. t. \quad - (K(B, C^T)\omega_1 + e_2 b_1) \geq e_2 - q_1 - d_1(1-\alpha) \\ q_{1i} \geq 0, \quad i = 1, \dots, m_2.$$

و مسأله ی دوگان بصورت (۴۳) در می آید.

$$Q(\omega_1, b_1, q_1, \beta_1, \gamma_1) = \frac{1}{2} \|K(A, C^T)\omega_1 + e_1 b_1\|^2 + C_1 e_2^T q_1 - \beta_1^T \{- (K(B, C^T)\omega_1 + e_2 b_1) + q_1 + d_1(1-\alpha) - e_2\} - \gamma_1^T q_1 \quad (43)$$

پس از مشتق گیری و برقراری شرایط مکمل زاید:

$$U_1 = -(S^T S + \epsilon I)^{-1} R^T \beta_1 \quad (44)$$

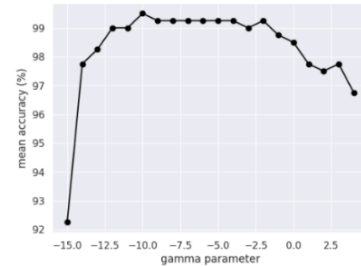
$$R = [K(B, C^T) e_2], S = [K(A, C^T) e_1] \quad \text{که در آن:}$$

## مراجع

- [۱] سبزه کار، مصطفی، بررسی تأثیر توجه مضاعف به نمونه‌های یادگیری با استفاده از قیود بهینه‌سازی طبقه‌بندی ماشین بردار پشتیبان، پایان‌نامه کارشناسی ارشد، گروه مهندسی کامپیوتر، دانشکده مهندسی، دانشگاه فردوسی مشهد، صفحات ۹۰-۸۲، آذر ۸۸
- [۲] ناجی عظیمی، زهرا، آشنایی با برنامه‌ریزی خطی فازی، ویراسته وحیدیان کامیاد، علی، ویرایش اول، مشهد، انتشارات دانشگاه فردوسی مشهد، تابستان ۹۵.
- [3] Alkenani, A. H., Li, Y., Xu, Y. & Zhang, Q. "Predicting Alzheimer's disease from spoken and written language using fusion-based stacked generalization", *J. Biomed. Inform.* 118, 103803, 2021.
- [4] Guo, Y. et al. "A review of wearable and unobtrusive sensing technologies for chronic disease management", *Comput. Biol. Med.* 129, 104163, 2020.
- [5] Higgins, V., Sohaei, D., Diamandis, E. P. & Prassas, I. "COVID-19: From an acute to chronic disease? Potential long-term health consequences", *Crit. Rev. Clin. Lab. Sci.* 58(5), 297-310, 2021.
- [6] Mozafari, Kourosh, Jalal A. Nasiri, Nasrollah Moghadam Charkari, and Saeed Jalili. "Action recognition by local space-time features and least square twin SVM (LS-TSVM)." In *2011 first international conference on informatics and computational intelligence*, pp. 287-292. IEEE, 2011.
- [7] Nasiri, Jalal A., and Amir M. Mir. "An enhanced KNN-based twin support vector machine with stable learning rules." *Neural computing and applications* 32, no. 16 : 12949-12969, 2020.
- [8] Raschka S, Mirjalili V. Python machine learning second edition, Packt Publishing, BIRMINGHAM – MUMBA, 2017.
- [9] Souza-Pereira, L., Pombo, N., Ouhbi, S., Felizardo, V. & Garcia, N. "Clinical decision support systems for chronic diseases: A systematic literature review", *Comput. Methods Progr. Biomed.* 195, 105565, 2020.
- [10] Tanveer M, Rajani T, Rastogi R, Shao YH, Ganaie MA. "Comprehensive review on twin support vector machines", *Annals of Operations Research*, 8:1-46, 2022.
- [11] Yuan, X., Chen, S., Yuwen, L., An, S., Mei, S. & Chen, T. "An improved SEIR model for reconstructing the dynamic transmission of COVID-19", In *Proceedings of IEEE International Conference on Bioinformatics and Biomedicine (BIBM)* 2320-2327, 2020.
- [12] Yuan X, Chen S, Sun C, Yuwen L. A novel early diagnostic framework for chronic diseases with class imbalance, *Scientific Reports*, 21;12(1):8614, 2022.

## پانویس ها

- Input space<sup>۱</sup>  
Kernel Trick<sup>۲</sup>  
Twin Support Vector Machine (TWSVM)<sup>۳</sup>  
Karush-Kuhn-Tucker (KKT)<sup>۴</sup>  
Relaxed Constraints Twin Support Vector Machine (RTWSVM)<sup>۵</sup>  
 $\alpha - cut$ <sup>۶</sup>  
Crisp<sup>۷</sup>



شکل ۴: اثر افزایش پارامتر  $\gamma$  روی دقت دسته‌بند *rtwsvm*

## ۳-۴- نتایج تجربی

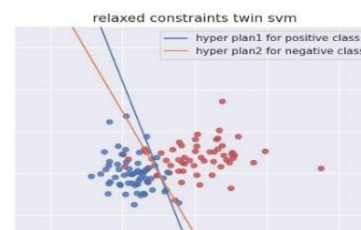
روش پیشنهادی با SVM ساده، رگرسیون خطی، درخت تصمیم و نزدیک‌ترین همسایه، که از مرجع [12] به دست آمده، مقایسه شده‌اند. نتایج مقایسه دقت و رتبه الگوریتم پیشنهادی در جدول شماره ۲ و ۳ قابل مشاهده‌اند؛ عملکرد بهتری دارد.

دقت	SVM	LR	DT	KNN	RTWSVM
CKD	۹۷.۸۵±۱.۵۵	۹۸.۷۵	۹۷.۵۰	۹۵.۰۰	۹۹.۵۰±۱
Hep	۸۰.۰۰	۸۳.۰۰	۸۳.۰۰	۸۰.۰۰	۱۰۰±۰
PIDD	۷۷.۹۲	۷۹.۸۷	۷۴.۶۸	۷۶.۶۲	۷۸.۸۱±۵.۵۹
WDBC	۹۷.۱۴	۹۷.۱۴	۹۵.۰۰	۹۷.۱۴	۹۸.۳۹±۱.۴۸
میانگین	۸۸.۲۲	۸۹.۶۹	۸۷.۵۴	۸۷.۱۹	۹۴.۱۷

جدول ۲: مقایسه دقت دسته‌بندی پیشنهادی و روش پیشنهادی

رتبه	SVM	LR	DT	KNN	RTWSVM
CKD	۳	۲	۴	۵	۱
Hep	۴	۲	۲	۴	۱
PIDD	۳	۱	۵	۴	۲
WDBC	۳	۳	۵	۳	۱
میانگین	۳.۳۷۵	۲.۱۲۵	۴.۱۲۵	۴.۱۲۵	۱.۲۵

جدول ۳: مقایسه رتبه دسته‌بندی پیشنهادی و روش پیشنهادی



شکل ۵: نمایش هندسی بردار پشتیبان دوقلو با قیود نرم.

## ۵- نتیجه‌گیری و پژوهش‌های آینده

گسترش فضای شدنی، عدم تأثیر منفی بر تابع هدف، کاهش اثر داده‌های پرت، سرعت بیشتر و کارایی بالاتر از ویژگی‌های مثبت روش جدید است. بکارگیری الگوریتم‌های جستجو مانند الگوریتم‌های تکاملی، می‌تواند یافتن مقدار بهینه برای پارامترهای زیاد موجود را سرعت بخشد. نیز می‌توان با در نظر گرفتن یک ماتریس وزنی ضرایب تحمل، به هریک از نمونه‌ها توجه ویژه داشت. با این کار تأثیر داده‌های نویزی کاهش خواهد یافت.