

دسته‌بند تک کلاسه مبتنی بر بردارهای پشتیبان برای داده‌های نویزی با استفاده از الگوریتم گروه میگوی آشوبی و تراکم محلی

محمدهادی قومنجانی^۱، دانشجوی کارشناسی ارشد؛ جواد حمیدزاده^۲، استادیار

۱- دانشکده مهندسی کامپیوتر و فناوری اطلاعات - دانشگاه صنعتی سجاد - مشهد - ایران - mghomanjani22@sadjad.ac.ir

۲- دانشکده مهندسی کامپیوتر و فناوری اطلاعات - دانشگاه صنعتی سجاد - مشهد - ایران - j_hamidzadeh@sadjad.ac.ir

چکیده: هدف دسته‌بندی تک کلاسه، تشخیص و جداسازی داده‌های اصلی از داده‌های پرت است. دسته‌بند توصیف داده‌ها مبتنی بر بردار پشتیبان، یکی از روش‌های دسته‌بندی تک کلاسه است. این روش با تعریف آبرگره‌ای در فضای ویژگی‌ها، سعی بر پوشش داده‌های اصلی در فضای آبرگره دارد. سطح آبرگره، مرز جداساز داده‌های اصلی از داده‌های پرت است. تعیین شعاع و مرکز مناسب برای آبرگره در فضای ویژگی‌ها، یک مسئله بهینه‌سازی است. وجود داده‌های نویزی در مجموعه داده‌های اصلی و عدم توجه به تراکم داده‌ها در انتخاب مرکز، از چالش‌های روش فوق است که باعث ایجاد خطا در تعیین مرز آبرگره خواهند شد. یکی از هدف‌های دسته‌بند پیشنهادی (KH-SVDD)، جستجوی مرکز مناسب برای آبرگره، با استفاده از الگوریتم بهینه‌سازی گروه میگوی آشوبی است. همچنین با استفاده از تراکم محلی نقاط داده‌ها، اهمیت و تأثیر نقاط بر مرز دسته‌بند، به صورت یک وزن محاسبه می‌شود. این وزن، پارامتری کمکی برای تشخیص داده‌های اصلی و نویزی است. برای ارزیابی روش پیشنهادی، آزمایش‌های متعددی بر روی مجموعه داده‌های واقعی انجام شده است. نتایج حاصل از آزمایش‌ها نشان‌دهنده عملکرد بهتر روش پیشنهادی از نظر تشخیص داده‌های نویزی در مقایسه با الگوریتم‌های مرز دانش است.

واژه‌های کلیدی: دسته‌بند تک کلاسه، تشخیص داده‌های پرت و نویزی، گروه میگوی آشوبی، وزن‌دهی تراکمی.

One-class Classifier Based on Support Vectors for Noisy Data by Using Chaotic Krill Herd Algorithm and Local Density

M. H. Ghomanjani¹, MSc Student; J. Hamidzadeh², Assistant Professor

1- Faculty of Computer Engineering and Information Technology, Sadjad University of Technology, Mashhad, Iran, Email: mghomanjani22@sadjad.ac.ir

2- Faculty of Computer Engineering and Information Technology, Sadjad University of Technology, Mashhad, Iran, Email: j_hamidzadeh@sadjad.ac.ir

Abstract: The purpose of one-class classification is to detect and separate target data from outlier. Support vector data description classifier is one of the one-class data classification methods. This method creates a hyper-sphere in feature space and tries to cover target data in the hyper-sphere. The hyper-sphere surface is the discernment boundary between target and outlier data. Determining appropriate radius and center for the sphere is an optimization problem. Existence of the noise in the data set and lack of attention to data density for choosing the center is the challenge of this method that triggered the mistake in determining of detection boundary. In the proposed classifier (KH-SVDD) we tried to search appropriate center of sphere with the use of chaotic krill herd optimization algorithm. Also, a weight is calculated for the effectiveness of the points on the classifier boundary with the use of local density of data points. This weight is an auxiliary parameter to detect target data from noise. The results of the experiments have been compared with state-of-the-art methods, which show superiority of the proposed method in noise detection.

Keywords: One-class classifier, outlier and noise detection, chaotic krill herd, weighted density.

تاریخ ارسال مقاله: ۱۳۹۵/۰۹/۱۷

تاریخ اصلاح مقاله: ۱۳۹۵/۱۱/۱۷ و ۱۳۹۶/۰۲/۱۰

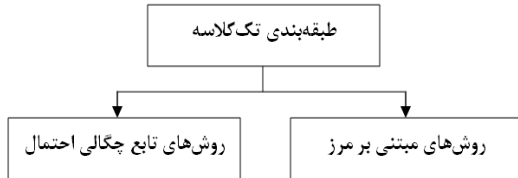
تاریخ پذیرش مقاله: ۱۳۹۶/۰۵/۰۹

نام نویسنده مسئول: جواد حمیدزاده

نشانی نویسنده مسئول: ایران - مشهد - بلوار جلال آل احمد ۶۴ - دانشگاه صنعتی سجاد - دانشکده مهندسی کامپیوتر و فناوری اطلاعات.

۱- مقدمه

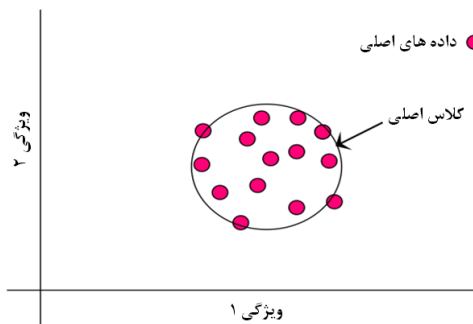
ساختار ادامه مقاله به شرح زیر است: در بخش دوم، مروری بر دسته‌بند توصیف داده‌ها مبتنی بر بردار پشتیبان، ارائه شده است. کارهای مرتبط پیشین در بخش سوم بررسی شده است. دسته‌بند پیشنهادی در بخش چهارم ارائه شده است. در بخش پنجم نتایج آزمایش‌ها نشان داده شده است. نتیجه‌گیری و کارهای آینده در بخش ششم مطرح گردیده است.



شکل ۱: تقسیم‌بندی روش‌های دسته‌بندی تک‌کلاسه

۲- دسته‌بند توصیف داده‌ها مبتنی بر بردار پشتیبان

دسته‌بند تک‌کلاسه توصیف داده‌ها مبتنی بر بردارهای پشتیبان (SVDD^۵)، از مهم‌ترین و پرکاربردترین روش‌ها جهت تشخیص داده‌های پرت است. این دسته‌بند جزء دسته دسته‌بندهای تک‌کلاسه مبتنی بر تشخیص مرز است. روش فوق به دنبال ساخت آبرگره‌ای برای پوشش داده‌های اصلی است. داده‌های درون آبرگره، به‌عنوان کلاس هدف و داده‌های خارج از آبرگره به‌عنوان داده‌های پرت، دسته‌بندی می‌شوند. مرز توصیف و تشخیص داده‌ها، همان رویه آبرگره است. این مرز با استفاده از نقاط داده‌ای که بر روی آن قرار دارند، توصیف می‌شود. به نقاط واقع بر روی مرز، بردارهای پشتیبان گفته می‌شود. شکل ۲ مثالی از آبرگره فوق، در فضای دوبعدی را نشان می‌دهد. داده‌های اصلی در شکل ۲ به‌صورت نقاط رنگی نشان داده شده‌اند. نقاط اصلی توسط یک دایره محاط شده‌اند. نقاطی که روی محیط دایره قرار دارند، بردارهای پشتیبان هستند و مرز دسته‌بند توسط این نقاط مشخص می‌شود [۱۰].



شکل ۲: آبرگره‌ای برای پوشش داده‌های اصلی

در این روش مرز آبرگره، با استفاده از لم هسته^۶ در فضای ویژگی‌های داده‌ها، قابل انعطاف می‌شود و امکان تفکیک پذیری خطی در فضاهای بالاتر (فضای ویژگی‌ها) افزایش می‌یابد. تفاوت مرز آبرگره،

تشخیص داده‌های پرت^۱ از مهم‌ترین فعالیت‌های داده‌کاوی محسوب می‌شود و در حوزه‌هایی مانند تشخیص بیماری، تشخیص نفوذ و اختلال سیستم، کاربردهای فراوانی دارد. یافتن داده‌های پرت در عمل کاری پرهزینه و در برخی موارد به‌علت در دسترس نبودن داده‌های پرت در مرحله آموزش، غیرممکن است. روش‌های متعددی از جمله مدل‌های آمار و احتمالی، مدل‌های خطی، مدل‌های تقریبی، مدل مارکوف، مدل K -نزدیک‌ترین همسایه و دسته‌بندی تک‌کلاسه برای تشخیص داده‌های پرت ارائه شده است. دسته‌بندی تک‌کلاسه^۲ یکی از پرکاربردترین روش‌ها برای تشخیص داده‌های پرت است [۴-۱].

دسته‌بندی^۳، انتساب یک نمونه داده جدید به یکی از کلاس‌های موجود است، این عملیات توسط روالی به نام دسته‌بند انجام می‌شود. به‌طور کلی در دسته‌بندی‌های تک‌کلاسه، مجموعه داده‌های در دسترس، داده‌های اصلی هستند و دانشی درباره داده‌های پرت وجود ندارد. در نتیجه داده‌های اصلی باید در یک کلاس (کلاس مثبت یا هدف) دسته‌بندی شوند. داده‌هایی که به کلاس تعیین‌شده، تعلق نداشته باشند به‌عنوان داده‌های پرت شناخته می‌شوند. از مهم‌ترین اهداف دسته‌بندی تک‌کلاسه، شناسایی و تشخیص مرز داده‌های اصلی از داده‌های پرت است [۵، ۶].

تفاوت دسته‌بندی تک‌کلاسه و دوکلاسه را می‌توان بدین‌صورت بیان کرد که در دسته‌بندی دوکلاسه (یا چندکلاسه)، ویژگی‌های دو کلاس اصلی (یا چند کلاس) در مرحله آموزش دسته‌بند، براساس مجموعه داده‌های برچسب‌دار تعیین می‌شود، سپس داده‌های دیگر با مقایسه با این دو کلاس (چند کلاس)، دسته‌بندی می‌شوند. در این نوع دسته‌بندی‌ها، همه داده‌ها باید در یکی از دو کلاس موجود، دسته‌بندی شوند، در نتیجه، در ماهیت دسته‌بندی دوکلاسه امکان تشخیص داده‌های پرت وجود ندارد [۷، ۸].

شکل ۱ تقسیم‌بندی کلی روش‌های دسته‌بندی تک‌کلاسه را نشان می‌دهد. اساس کار دسته‌بندها در روش‌های مبتنی بر مرز، به‌دست‌آوردن مرزی برای جداسازی داده‌های کلاس اصلی از داده‌های پرت است. گروه دوم، دسته‌بندهای مبتنی بر احتمال است. در این دسته‌بندها، با استفاده از تابع چگالی احتمال، یک مقدار عددی برای داده‌های موجود، محاسبه می‌شود. با مقایسه این مقدار عددی با یک مقدار آستانه، تعلق و عدم تعلق نمونه داده‌ها به کلاس اصلی مشخص می‌شود [۹].

دسته‌بند توصیف داده‌ها مبتنی بر بردارهای پشتیبان، یکی از دسته‌بندهای تک‌کلاسه است، که مهم‌ترین کاربرد آن تشخیص مرز بین داده‌های اصلی و داده‌های پرت است. در این مقاله یک دسته‌بند تک‌کلاسه (KH-SVDD^۴) بر پایه دسته‌بند توصیف داده‌ها مبتنی بر بردارهای پشتیبان ارائه شده است. ایده اصلی این دسته‌بند، تشخیص درست داده‌های نویزی است.

پشتیبان) برای مرز آبرگره پیدا شود. تابع ϕ ، یک تابع نگاشت است که یک نقطه را به ابعاد بالاتر در فضای ویژگی نگاشت می‌کند.

با توجه به معادله (۱)، هرگاه فاصله یک نقطه تا مرکز آبرگره، کمتر مساوی شعاع آبرگره باشد، آن نقطه به کلاس اصلی تعلق دارد. به عبارت دیگر، نقاطی که فاصله ویژگی‌های آنها تا مرکز آبرگره، بزرگ‌تر از شعاع کره باشند، به‌عنوان نقاط پرت شناخته می‌شوند. با حل مسئله بهینه‌سازی معادله (۱) با استفاده از ضرایب لاگرانژ، معادله (۱) به فرم مشتق‌پذیر معادله (۲) تبدیل می‌شود [۱۳].

$$L = R^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [R^2 + \xi_i - (\phi(x_i) - a)^2] - \sum_{i=1}^N \gamma_i \xi_i, \alpha_i, \gamma_i \geq 0 \quad (2)$$

در این معادله α_i, γ_i ضرایب لاگرانژ هستند. هدف معادله (۲)، کمینه‌سازی عبارت L با توجه به کمینه‌سازی متغیرهای R, a, ξ_i و بیشینه‌سازی ضرایب لاگرانژ است. با مشتق‌گیری جزئی معادله (۲) نسبت به متغیرهای R, a, ξ_i ، مقدار شعاع، موقعیت مرکز و مقدار خطای نقاط خارج از آبرگره حاصل می‌شود.

$$\frac{\partial L}{\partial R} = 0 \rightarrow \sum \alpha_i = 1 \quad (3)$$

معادله (۳) مشتق جزئی نسبت به متغیر شعاع است، در نتیجه مجموع ضرایب لاگرانژ در مسئله یک می‌شود.

$$\frac{\partial L}{\partial a} = 0 \rightarrow a = \frac{\sum \alpha_i \phi(x_i)}{\sum \alpha_i} = \sum \alpha_i \phi(x_i) \quad (4)$$

معادله (۴) مشتق جزئی نسبت به مرکز آبرگره است که با کمک آن موقعیت مرکز آبرگره محاسبه می‌شود.

$$\frac{\partial L}{\partial \xi_i} = 0 \rightarrow \alpha_i + \gamma_i = C \rightarrow 0 \leq \alpha_i \leq C, \gamma_i, C \geq 0 \quad (5)$$

با توجه به معادله (۵) محدودیتی برای ضریب لاگرانژ α تعیین می‌شود. بعد از جایگزینی معادله‌های فوق در معادله (۲)، دوگان مسئله به صورت معادله (۶) ایجاد می‌شود.

$$\max L = \sum_{i=1}^N \alpha_i (\phi(x_i) \cdot \phi(x_j)) - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j (\phi(x_i) \cdot \phi(x_j)) \quad (6)$$

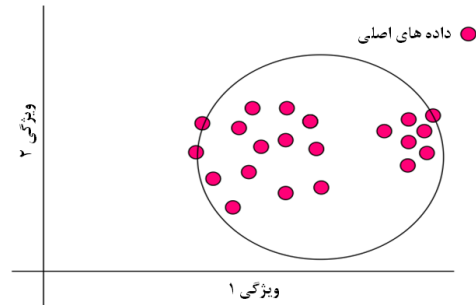
$$\sum \alpha_i = 1, \quad 0 \leq \alpha_i \leq C$$

در معادله (۶)، $\phi(x_i) \cdot \phi(x_j)$ به صورت یک تابع هسته تعریف می‌شود. معادله (۶) یک مسئله بهینه‌سازی محدب درجه دو است که یک جواب بهینه سراسری دارد. ارزیابی نقاط توسط این دسته‌بندی، با استفاده از معادله (۷) صورت می‌گیرد. به‌ازای هر نقطه، اگر $f(x) \leq 0$ شود، آن نقطه به کلاس هدف تعلق دارد در غیر این صورت به‌عنوان نقطه پرت معرفی می‌شود.

$$f(x) = \|\phi(x_i) - a\|^2 - R^2 \quad (7)$$

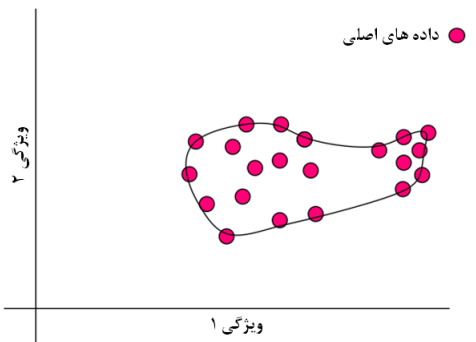
بدون و با استفاده از لم هسته در شکل‌های ۳ و ۴ نشان داده شده است [۱۱].

شکل ۳ آبرگره‌ای محاط‌شده بر مجموعه داده‌های اصلی، بدون استفاده لم هسته را نشان می‌دهد. در این حالت امکان تفکیک‌پذیری و انعطاف مرز آبرگره، در فضای ویژگی‌ها وجود ندارد. شکل ۴ همان آبرگره، با استفاده لم هسته را نشان می‌دهد. مرز آبرگره در شکل ۴ منعطف شده و امکان تفکیک‌پذیری در فضای ویژگی‌ها را دارا است [۱۲].



شکل ۳: مرز آبرگره بدون استفاده از لم هسته

روش SVDD به دنبال پوشش حداکثری داده‌های اصلی و کاهش شانس پذیرش داده‌های پرت در کلاس هدف است. در SVDD، یک مسئله بهینه‌سازی به صورت معادله (۱) تعریف می‌شود [۱۳].



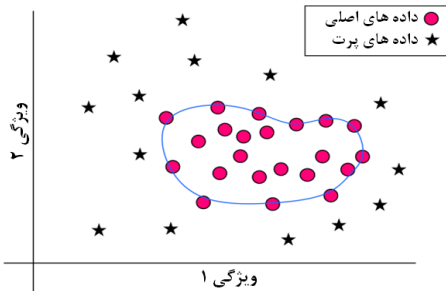
شکل ۴: مرز آبرگره با استفاده از لم هسته

$$\min R_2 + C \sum_{i=1}^N \xi_i \quad (1)$$

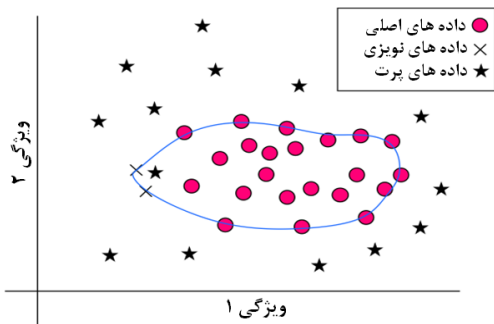
$$s.t. \quad \|\phi(x_i) - a\|^2 \leq R^2 + \xi_i, \quad \xi_i \geq 0, \quad i = 1, 2, \dots, N$$

هدف مسئله بهینه‌سازی فوق، تعریف آبرگره‌ای با حداقل شعاع R به مرکز a برای احاطه کردن داده‌های اصلی در فضای ویژگی‌ها است. داده‌هایی که خارج از آبرگره قرار گیرند، خطایی به مقدار ξ_i دارند. مقدار ξ_i در واقع فاصله نقاط بیرون از آبرگره تا مرکز آبرگره است. نقاط بیرون از آبرگره با ضریبی به مقدار C ، جریمه می‌شوند و در مسئله بهینه‌سازی استفاده می‌شوند تا نقاط مرزی مناسبی (بردارهای

در دسته‌بند پیشنهادی با آگاهی از این اصل که داده‌های نویزی از تراکم کمتری نسبت به داده‌های اصلی در مجموعه داده‌ها برخوردار هستند [۱۷]، از معیارهای تراکم محلی نقاط و فاصله تا مرکز آبرگره استفاده شده، تا علاوه بر تشخیص داده‌های پرت توسط دسته‌بند، نقاط نویزی در مجموعه داده‌های اصلی شناسایی شوند و از تأثیر آنها بر مرز دسته‌بند، جلوگیری شود.



شکل ۶: مرز آبرگره بدون وجود نویز در مجموعه داده‌های اصلی



شکل ۷: تأثیر نویز موجود در مجموعه داده‌ها بر مرز آبرگره

۳- مرور و بررسی کارهای مرتبط پیشین

تلاش‌های زیادی جهت بهبود روش توصیف داده‌ها به کمک بردارهای پشتیبان صورت گرفته است. در سال ۲۰۰۴ تکس و دوین در روشی به نام NSVDD، با فرض این‌که داده‌های پرت نیز در دسترس می‌باشند، از آنها در فرآیند یادگیری استفاده کردند. هدف این روش ساخت آبرگره‌ای برای پوشش داده‌های اصلی و عدم پوشش داده‌های پرت است [۱۳].

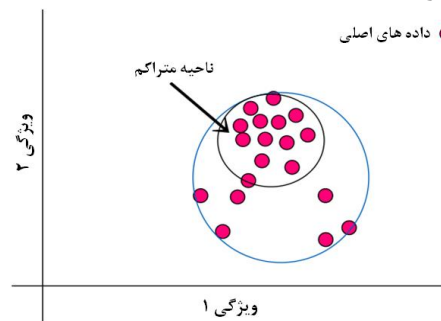
در سال ۲۰۰۷ لی و همکاران [۱۸] یک روش جدید بر پایه SVDD معرفی کردند، در این الگوریتم رویکرد جدیدی برای معیار فاصله براساس نزدیک‌ترین همسایه و پنجره پارزن^۷ و درجه تراکم تمام نقاط داده‌ها، ارائه شد.

در سال ۲۰۰۹ وانگ و همکاران با تولید داده‌های پرت بیرون از مرز آبرگره، به صورت مصنوعی و دخالت آنها در مرحله آموزش دسته‌بند توصیف بردارهای پشتیبان، عملکرد آن را بهبود دادند [۱۹].

هم‌چنین در سال ۲۰۰۹ در [۲۰]، با ترکیب فیلترینگ متعامد و روش SVDD، سعی در بهبود عملکرد دسته‌بند در برخورد با داده‌های

با توجه به معادله فوق نقاطی که ضریب لاگرانژ α_i آنها صفر باشد، به‌عنوان نقاط داخل آبرگره و نقاطی که ضریب لاگرانژ آنها بین صفر و C باشد، نقاط مرزی آبرگره هستند و به‌عنوان بردارهای پشتیبان شناخته می‌شوند. هم‌چنین اگر ضریب لاگرانژ برابر C باشد، نقاط خارج از آبرگره قرار دارند و به‌عنوان نقاط پرت شناخته می‌شوند.

در این مقاله یک دسته‌بند تک‌کلاسه برای تشخیص داده‌های نویزی (KH-SVDD) ارائه شده است. یکی از نقاط ضعف دسته‌بند SVDD عدم توجه به تراکم نقاط در مجموعه داده‌ها، برای انتخاب مرکز آبرگره است. یک ناحیه متراکم نزدیک به مرز آبرگره باعث خطا در انتخاب مرز، در نتیجه افزایش خطا در تشخیص داده‌های اصلی از پرت می‌شود [۲]، [۳]. عدم توجه به تراکم در انتخاب مرکز آبرگره در شکل ۵ نشان داده شده است. در این شکل مرز آبرگره به ناحیه پرتراکم داده‌ها نزدیک است.



شکل ۵: ناحیه متراکم داده‌ها نزدیک به مرز آبرگره

در دسته‌بند پیشنهادی، از ایده الگوریتم غذایابی میگوی آشوبی برای یافتن مرکز بهینه آبرگره استفاده شده است. در طبیعت، میگوها در اطراف یک منبع غذایی یک ناحیه متراکم ایجاد می‌کنند. یک میگو با جستجو و انتخاب یک ناحیه متراکم میگوها، به سمت مرکز تراکم (محل غذا) حرکت می‌کند. از ایده فوق برای یافتن مرکز متراکم و مناسب برای آبرگره استفاده شده است.

وجود نویز در مجموعه داده‌ها، یکی از چالش‌های مهم در دسته‌بندی است. عوامل مختلفی از جمله خطا در دستگاه‌های اندازه‌گیری، نیروی انسانی و نویزهای محیطی مانند دما، فشار و گرما موجب ایجاد خطا و نویز در مجموعه داده‌های اصلی جمع‌آوری شده، می‌شوند [۱۴].

حضور داده‌های نویزی در مجموعه داده‌های اصلی، چالشی بسیار مهم برای دسته‌بندهای تک‌کلاسه است. رفتار دسته‌بند با داده‌های نویزی مانند داده‌های اصلی است. در نتیجه، با توجه به تأثیر داده‌های نویزی در مرحله آموزش دسته‌بند، مرز دسته‌بند دچار خطا می‌شود (شکل ۷). شکل ۶ مرز دسته‌بند بدون حضور داده‌های نویزی را نشان می‌دهد، اما در شکل ۷ داده‌های نویزی موجب خطا، در تشخیص مرز دسته‌بند شده‌اند. همان‌طور که در شکل ۷ مشخص است یک نقطه داده پرت داخل آبرگره کلاس اصلی است که نشان از خطای دسته‌بند است [۱۵]، [۱۶].

۴- روش پیشنهادی

دسته‌بند پیشنهادی (KH-SVDD) از دو گام اصلی تشکیل شده است. در گام اول، هدف پیدا کردن یک مرکز بهینه برای آبرگره با توجه به تراکم نقاط است و هدف گام دوم، ساخت مرز آبرگره با کمترین شعاع ممکن و بیشترین دقت است. هم‌چنین مرحله آموزش دسته‌بند پیشنهادی در این دو گام، با فرض در دسترس بودن داده‌های اصلی انجام می‌شود. روش پیشنهادی در دو بخش، براساس گام‌های اصلی بررسی می‌شود. بخش ۴-۱ به بررسی، چگونگی جستجوی مرکز برای ساخت آبرگره می‌پردازد. در بخش ۴-۲ نحوه تعیین مرز آبرگره مشخص می‌شود.

۴-۱- جستجوی مرکز آبرگره

هدف از گام اول، جستجوی مرکزی مناسب برای آبرگره است. مرکز یک ناحیه متراکم، می‌تواند به‌عنوان یک جواب بهینه برای مرکز آبرگره باشد. جستجوی نقطه‌ای متراکم، در مسئله با ابعاد بالا (تعداد ویژگی‌های زیاد) و تعداد نقاط زیاد، یک کار زمان‌بر است. در نتیجه استفاده از روش‌هایی که کل مجموعه داده‌ها را برای شمارش تراکم نقاط، پیمایش کنند یک کار زمان‌بر است. در روش پیشنهادی از الگوریتم گروه میگوی آشوبی برای یافتن نقاط متراکم محلی استفاده شده است.

در این مقاله از الگوی رفتار میگوها (در جستجوی غذا و پیوستن به یک دسته میگو) در جهت پیدا کردن یک مرکز آبرگره بهینه، برای توصیف داده‌ها الهام گرفته شده است [۲۷]. مسئله دسته‌بندی تک‌کلاسه مانند یک مسئله غذاییابی گروه میگوها در نظر گرفته شده است.

یکی از مشکلات الگوریتم SVDD، عدم توجه به تراکم نقاط در انتخاب مرکز آبرگره است. در دسته‌بند پیشنهادی (KH-SVDD) تلاش شده است تا مرکز آبرگره به سمت ناحیه پرتراکم حرکت کند. این ایده موجب می‌شود مرز آبرگره از نواحی متراکم دور شود که موجب کاهش خطا در تشخیص نقاط اصلی از نقاط پرت و نویزی می‌شود. در دسته‌بندی‌های تک‌کلاسه که فقط داده‌های یک کلاس موجود هستند، تراکم نقاط داده‌های اصلی، بسیار بیشتر از تراکم نقاط نویزی است.

هر نقطه داده (مانند یک میگو که به یک دسته میگو برای غذاییابی ملحق می‌شود) در جستجوی یک مرکز آبرگره، برای الحاق به آن است. رفتار غذاییابی یک میگو (که در جستجوی غذا در ناحیه پرتراکم است و به سمت مرکز ناحیه متراکم حرکت می‌کند) مانند رفتار مرکز آبرگره (مرکز اولیه با استفاده از الگوریتم SVDD مقداردهی می‌شود) است که تمایل دارد نقش خود را با یک نقطه متراکم‌تر عوض کند.

در الگوریتم گروه میگوی آشوبی، تابع ارزیابی تغییرات حرکت میگو، براساس فاصله میگو تا غذا (F)، تراکم میگوها (D) و حرکت ایجادشده توسط سایر میگوها (N) به‌صورت معادله (۹) در فضای n -بعدی تعریف می‌شود.

نویزی شد. فیلترینگ متعامد به‌عنوان یک پیش‌پردازش در این روش است. این فیلتر انحراف معیار ناخواسته در داده‌ها را حذف می‌کند.

در روشی دیگر [۲۱] با یک روش پس‌پردازشی^۸ سعی در بهبود عملکرد SVDD شد. در این روش سعی شده است بار محاسباتی لم هسته در منعطف‌سازی مرز آبرگره حذف شود. در این الگوریتم برای اصلاح مرزهای آبرگره از الگوریتم نزدیکترین همسایه استفاده شده است.

در تحقیق [۲۲]، روشی براساس SVDD، برای توصیف داده‌های دو کلاس معرفی شد. در این روش برای هر کلاس به‌طور جداگانه، آبرگره‌ای در نظر گرفته می‌شود و داده‌های خارج از این دو آبرگره به‌عنوان داده‌های پرت تلقی می‌شوند.

در سال ۲۰۱۳ لیو و همکاران یک روش جدید بر پایه SVDD معرفی کردند. در این روش برای نقاط داده‌های اصلی یک امتیاز اعتبار، تعریف می‌شود. این امتیاز، احتمال درست‌نمایی^۹ نقاط برای عضویت در کلاس هدف را مشخص می‌کند [۲۳].

چا و همکاران در سال ۲۰۱۴ روشی به نام DW-SVDD را براساس الگوریتم توصیف داده‌ها مبتنی بر بردارهای پشتیبان، ارائه کردند [۲۴]. در این روش نگاهی به توصیف نقاط داده‌ها براساس چگالی آنها شده است. رویکرد استفاده‌شده در این روش برای توزیع تراکم داده‌ها، استفاده از روش KNN است. با توجه به تراکم نقاط به آنها وزنی مشخص داده می‌شود تا اولویت تأثیر آنها بر مرز دسته‌بند مشخص شود. در این روش برای انتخاب مرکز آبرگره از تراکم داده‌ها استفاده نمی‌شود.

در سال ۲۰۱۵ چن و همکاران روشی به نام R-SVDD را بر پایه الگوریتم توصیف داده‌ها مبتنی بر بردار پشتیبان، برای دسته‌بندی در مجموعه داده‌های نویزی (هنگامی که فقط مجموعه داده‌های اصلی در دسترس است) معرفی کردند [۲۵]. در این روش از معیار فاصله‌ای به نام Cut-off استفاده شده است و با کمک همین معیار فاصله، تراکم محلی نقاط محاسبه شده و سپس از آن در تابعی برای جریمه نقاط داده‌ها استفاده می‌شود. از اشکالات این روش عدم توجه به تراکم داده‌ها در انتخاب مرکز آبرگره است. هم‌چنین برای تشخیص داده‌های نویزی از داده‌های اصلی از تابع معادله (۸) برای جریمه نقاط استفاده می‌شود. در این تابع توجهی به فاصله نقاط تا مرکز آبرگره نشده است و فقط از تراکم محلی نقاط ($\rho(x_i)$) استفاده شده است.

$$w(x_i) = \rho(x_i) / \max(\rho(x_i)) \quad (8)$$

کیم و همکاران در سال ۲۰۱۵ روشی جدیدی با استفاده از یادگیری عمیق برای توصیف داده‌ها معرفی کردند [۲۶]. در این روش برای غلبه بر مشکل بیش‌برازش^{۱۰} از یک مدل چندلایه شبکه عصبی با یادگیری عمیق استفاده شده است. در این مدل لایه‌های پنهان وظیفه استخراج ویژگی‌های لازم برای SVDD را بر عهده دارند. هم‌چنین هر لایه به‌عنوان دسته‌بند یک کلاس، برای مجموعه داده‌ها است.

$$x_i(t + \Delta t) = x_i(t) + \Delta t \cdot \frac{dx_i}{dt} \quad (14)$$

معادله حرکت یک میگو در طول زمان با کمک معادله (۱۴) بیان می‌شود، در روش پیشنهادی معادله جابه‌جایی مرکز آبرگره در تکرارهای مختلف (براساس معادله (۱۴))، با استفاده از معادله (۱۵) نشان داده می‌شود، که در آن t تعداد دفعات اجرای الگوریتم، را نشان می‌دهد.

$$x_i(t + 1) = x_i(t) + dx_i \quad (15)$$

به‌طور کلی در گام اول روش پیشنهادی، یک بردار ویژگی‌ها، برای مرکز اولیه آبرگره (مانند یک کروموزوم) مقداردهی می‌شود. این بردار توسط تابع برازش (۱۰) مورد ارزیابی قرار می‌گیرد. در یک گام اجرای الگوریتم تراکم یک نقطه (مرکز فعلی) با نقاط همسایه‌اش مقایسه می‌شود و نقطه‌ای که بیشترین مقدار تابع برازش را به‌دست آورد به‌عنوان نسل برگزیده، مرکز آبرگره را تغییر می‌دهد. هم‌چنین در روش پیشنهادی از عملگرهای تولیدمثل، باز ترکیب و جهش معرفی شده در الگوریتم گروه میگوی آشوبی [۲۲]، استفاده می‌شود.

یکی از مشکلات روش‌های جستجو، گرفتار شدن در نقاط متراکم محلی است. در روش پیشنهادی برای جلوگیری از انتخاب نقطه متراکم محلی به‌جای سراسری، از توابع نگاشت آشوب برای تولید تصادفی گام جستجو، استفاده می‌شود. توابع آشوب این امکان را فراهم می‌سازند تا با در اختیار داشتن تابع نگاشت آشوب و مقدار اولیه تابع، اعداد تصادفی دوباره تولید شوند. این قابلیت برای جلوگیری از تولید نقاط تصادفی تکراری در الگوریتم که باعث انتخاب تکراری یک نقطه متراکم شود، بسیار مفید است.

از تابع آشوب معادله (۱۶) (Logistic map[20]) برای تغییر موقعیت تصادفی مرکز آبرگره، استفاده می‌شود. ($r = 3/9999$)

$$X_{n+1} = r \cdot X_n \cdot (1 - X_n) \quad (16)$$

۴-۲- تعیین مرکز آبرگره (تعیین شعاع)

هدف این گام، یافتن یک شعاع کمینه برای تشکیل مرکز آبرگره است. در ابتدا یک مقدار جریمه برای نقاط، علاوه بر میزان خطای ξ_i ، تعریف می‌شود. این مقدار براساس تراکم محلی نقاط و فاصله تا مرکز آبرگره با توجه به معادله (۱۷) محاسبه می‌شود. هر چقدر میزان این جریمه کمتر باشد، نقطه مورد نظر تأثیری کمتری بر مرکز دسته‌بند در مرحله آموزش دارد.

$$\omega(x_i) = \frac{D_i}{\max(D_i)} \cdot \frac{1}{\text{dist}(x_i)} \quad (17)$$

در معادله (۱۷)، D_i تراکم محلی یک نقطه و $\text{dist}(x_i)$ فاصله نقطه مورد نظر تا مرکز آبرگره است. مقدار جریمه محاسبه‌شده، عددی بین صفر و یک است. با استفاده از جریمه محاسبه‌شده معادله (۱) را به‌صورت معادله (۱۸) تغییر می‌دهیم.

$$\frac{dx_i}{dt} = N_i + F_i + D_i \quad (9)$$

در دسته‌بند پیشنهادی، تغییرات زمان (dt) وجود ندارد و مانند سایر الگوریتم‌های تکاملی، جایگزین تغییرات زمان، تغییر در تکرار اجراهای الگوریتم است. در نتیجه تغییرات زمان (اختلاف تکرارهای الگوریتم) یک واحد است ($dt = 1$).

ارزیابی تغییرات حرکت نقاط مرکز، در روش پیشنهادی توسط معادله (۱۰) بر پایه معادله (۹) محاسبه می‌شود، که در آن فاصله یک نقطه داده از مرکز قبلی (F_i)، تراکم محلی یک نقطه داده (D_i) و حرکت ایجادشده سایر نقاط (N_i) است. حرکت سایر نقاط در واقع جابه‌جایی نقش مرکز آبرگره، توسط نقاط است، که به سمت مرکز ناحیه متراکم جابه‌جا می‌شوند.

$$dx_i = N_i + F_i + D_i \quad (10)$$

با توجه به این که در روش پیشنهادی نقاط با تراکم بالا برای مسئله دارای اهمیت می‌باشند، تراکم محلی نقاط داده، با استفاده از نامساوی (۱۱) محاسبه می‌شود. هر دو نقطه از داده‌ها که در نامساوی (۱۱) صدق کنند به‌عنوان همسایه یکدیگر شناخته می‌شوند. به تعداد نقاطی که در همسایگی یک نقطه باشند، تراکم آن نقطه گفته می‌شود.

$$\frac{1}{d} \sum_{f=1}^d \|x_{if} - x_{jf}\| \leq \eta \quad (11)$$

در نامساوی (۱۱)، d تعداد ویژگی‌های داده‌ها و ویژگی f -ام از نقطه i -ام است. در معادله فوق ابتدا میانگین فاصله بین تمام ویژگی‌های دو نقطه از هم محاسبه می‌شود، اگر مقدار محاسبه‌شده از یک مقدار آستانه (η) کمتر یا مساوی باشد، دو نقطه در همسایگی یکدیگر قرار دارند. هم‌چنین از معادله سمت چپ نامساوی (۱۱) برای محاسبه فاصله یک نقطه تا مرکز آبرگره نیز استفاده می‌شود.

در الگوریتم گروه میگو، حرکت میگوها به سمت مرکز تراکم (محل غذا) است. معادله این حرکت با استفاده از تراکم محلی نقطه فعلی و تراکم محلی نقطه هدف، با کمک معادله (۱۲) تقریب زده می‌شود.

$$N_i^{new} = \alpha_i \cdot N^{max} \cdot \omega \cdot N^{old} \quad (12)$$

در معادله (۱۲)، ω یک وزن تصادفی در بازه صفر و یک است. N^{max} بیشترین سرعت و N^{old} آخرین حرکت انجام شده است. جهت حرکت (α_i)، به کمک معادله (۱۳) محاسبه می‌شود.

$$\alpha_i = \alpha_i^{local} + \alpha_i^{target} \quad (13)$$

در روش پیشنهادی صورت معادله‌های (۱۲) و (۱۳) تغییر نمی‌کند و پارامترهای آن به‌این‌صورت می‌باشند: N^{max} بیشترین فاصله بین دو مرکز آبرگره قبلی و جدید، در تمام تکرارهای الگوریتم است. N^{old} فاصله بین مرکز قبلی و جدید، در تکرار اخیر الگوریتم است. هم‌چنین پارامتر α_i معادله (۱۳)، برای روش پیشنهادی، جهت حرکت نقاط مرکز (جهت تغییر نقش نقاط برای مرکز شدن در فضای ویژگی‌های داده‌ها) را نشان می‌دهد.

شده است. در این حالت یک کلاس از مجموعه داده‌ها (کلاس ۱ یا کلاس ۲ مشخص شده در جدول ۱) به عنوان کلاس هدف در مرحله آموزش دسته‌بند استفاده می‌شود و کلاس دیگر به عنوان کلاس پرت در نظر گرفته می‌شود و در مرحله آزمایش دسته‌بند استفاده می‌شود. برای افراز مجموعه داده‌ها به مجموعه آموزشی و آزمایشی از روش اعتبارسنجی متقاطع ۱۰-تایی^{۱۱} استفاده شده است. برای مقایسه عملکرد دسته‌بندها از معیار صحت دسته‌بند (معادله (۲۲)) و معیار نرخ مثبت صحیح ((معادله (۲۳)) استفاده شده است [۲۹].

جدول ۱: مجموعه داده‌های استفاده‌شده در آزمایش‌ها برگرفته از UCI

شماره	مجموعه داده	تعداد ویژگی	کلاس ۱		کلاس ۲	
			نام دسته	تعداد نمونه	نام دسته	تعداد نمونه
۱	Breast Cancer	۳۲	Benign	۳۵۷	Malignant	۲۱۲
۲	Ionosphere	۳۴	Good	۲۲۵	Bad	۱۲۶
۳	Pima Indians	۸	Yes	۵۰۰	No	۲۶۸
۴	Spambase	۵۷	Non-spam	۲۷۸۸	Spam	۱۸۱۳
۵	Sonar	۶۰	Rock	۹۷	Metal	۱۱۱

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

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

تعریف اصطلاحات مثبت صحیح (TP)، مثبت کاذب (FP)، منفی صحیح (TN) و منفی کاذب (FN) براساس جدول ۲ است.

جدول ۲: ماتریس درهم‌ریختگی^{۱۲} نقاط پیش‌بینی شده

داده‌های پرت	داده‌های پرت	داده‌های اصلی
پیش‌بینی شده	پیش‌بینی شده	پیش‌بینی شده
داده‌های پرت واقعی	مثبت صحیح	منفی کاذب
داده‌های اصلی واقعی	مثبت کاذب	منفی صحیح

آزمایش‌ها بر روی دسته‌بندها در دو فرضیه صورت می‌پذیرد. در فرضیه اول مجموعه داده‌های اصلی فاقد نویز می‌باشند و در فرضیه دوم در مجموعه داده‌های اصلی، ۱۰ درصد از نقاط آن را به صورت تصادفی تغییر می‌دهیم تا مجموعه داده‌ها شامل نویز شود (با استفاده از نویز توزیع گوسی بر روی ویژگی‌ها [۳۰]). جدول ۳ و ۴ نتایج آزمایش‌ها بر روی مجموعه داده‌ها برای فرضیه اول (داده‌های بدون نویز) را نشان می‌دهد (براساس آموزش کلاس ۱ به عنوان کلاس اصلی). در جدول ۳ میزان صحت دسته‌بندها در مجموعه داده‌ها ارائه شده است.

$$\min R_2 + C \sum_{i=1}^N \omega(x_i) \xi_i \quad (18)$$

$$s.t. \quad \|\phi(x_i) - a\|^2 \leq R^2 + \xi_i, \quad \xi_i \geq 0, \quad i = 1, 2, \dots, N$$

با اضافه کردن ضرایب لاگرانژ به معادله (۱۸)، معادله بهینه‌سازی (۱۹) حاصل می‌شود.

$$L = R^2 + C \sum_{i=1}^N \omega(x_i) \xi_i - \sum_{i=1}^N \alpha_i [R^2 + \xi_i - (\phi(x_i) - a)^2] - \sum_{i=1}^N \gamma_i \xi_i, \quad \alpha_i, \gamma_i \geq 0 \quad (19)$$

با مشتق‌گیری جزئی از متغیرهای شعاع، مرکز و خطا معادله (۱۹)، دوگان مسئله به صورت معادله (۲۰) تبدیل می‌شود.

$$\max L = \sum_{i=1}^N \alpha_i (\phi(x_i) \cdot \phi(x_j)) - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j (\phi(x_i) \cdot \phi(x_j)) \quad (20)$$

$$\sum \alpha_i = 1, \quad 0 \leq \alpha_i \leq \omega(x_i) C$$

شعاع آبرگره با توجه به بردارهای پشتیبان (SVs) با کمک معادله (۲۱) محاسبه می‌شود.

$$R = \sqrt{\frac{1}{|SVs|} \sum \|\phi(x_i) - a\|^2} \quad (21)$$

گام‌های دسته‌بند KH-SVDD در الگوریتم ۱ نشان داده شده است.

الگوریتم ۱: الگوریتم دسته‌بند KH-SVDD

- At the first epoch: calculate center of the hyper-sphere (with SVDD Eq. (4))
- First Step:
- while (not Stop-Condition)
- for each data point is neighbor with center of the hyper-sphere
- Calculate density of data points by Eq. (11)
- Select maximum local density as new center
- Evaluate fitness function for center by Eq. (10)
- end for
- use Eq. (16) for generate random center of hyper-sphere for find dense point
- end while
- Second Step:
- for each data point calculate Eq. (17)
- Determine SVs by Eq. (7)
- Calculate R by Eq. (21)
- end for

۵- ارزیابی و مقایسه روش پیشنهادی

در این بخش نتایج حاصل از آزمایش‌ها بر روی دسته‌بند KH-SVDD با روش‌های [13] SVDD, [24] DW-SVDD, [25] R-SVDD آورده شده است. مجموعه داده‌های جدول ۱ برگرفته از UCI به عنوان مجموعه داده‌های واقعی، برای آزمایش‌ها استفاده شده است [۲۸]. مجموعه داده‌ها در دو کلاس، شماره ۱ و ۲ دسته‌بندی شده‌اند که در زمان آموزش فقط از یک از این دو کلاس برای آموزش دسته‌بند استفاده

جدول ۵: میانگین صحت دسته‌بندها در فرضیه داده‌های بدون نویز (کلاس ۲-
کلاس اصلی)

KH-SVDD	R-SVDD[25]	DW-SVDD[24]	SVDD[13]	مجموعه داده‌ها
۰/۹۷۹ (۱)	۰/۹۷۱ (۳)	۰/۹۷۴ (۲)	۰/۹۴۹ (۴)	Breast Cancer
۰/۹۸۶ (۱)	۰/۹۷۹ (۳)	۰/۹۸۵ (۲)	۰/۹۴۲ (۴)	Ionosphere
۰/۹۵۹ (۲)	۰/۹۵۹ (۱)	۰/۹۵۳ (۳)	۰/۹۳۱ (۴)	Pima Indians Diabetes
۰/۹۶۵ (۲)	۰/۹۶۳ (۳)	۰/۹۶۶ (۱)	۰/۹۴۳ (۴)	Spambase
۰/۹۶۰ (۱)	۰/۹۵۴ (۳)	۰/۹۵۷ (۲)	۰/۹۵۱ (۴)	Sonar

همان طور که در جدول‌های ۱۰-۷ مشخص است صحت و نرخ مثبت صحیح دسته‌بندها در دسته‌بندی داده‌های دارای نویز به‌میزان قابل توجهی کاهش یافته است، دسته‌بند R-SVDD و دسته‌بند پیشنهادی KH-SVDD به‌دلیل داشتن رویکردی برای تشخیص داده‌های نویزی، نتایج بهتری را در مواجهه با مجموعه داده‌های دارای نویز، از خود نشان داده‌اند. هم‌چنین نتایج دسته‌بند پیشنهادی به‌میزان قابل قبولی از دیگر دسته‌بندها بهتر است. در آزمایش دیگر، میزان نویز در مجموعه داده‌ها به‌صورت پله‌ای افزایش یافته است، تا تأثیر افزایش نویز بر رفتار دسته‌بندها بررسی شود.

جدول ۶: میانگین نرخ مثبت صحیح دسته‌بندها در فرضیه داده‌های بدون نویز (کلاس ۲- کلاس اصلی)

KH-SVDD	R-SVDD[25]	DW-SVDD[24]	SVDD[13]	مجموعه داده‌ها
۰/۹۸۳ (۱)	۰/۹۷۹ (۳)	۰/۹۸۱ (۲)	۰/۹۶۲ (۴)	Breast Cancer
۰/۹۸۵ (۲)	۰/۹۸۲ (۳)	۰/۹۸۹ (۱)	۰/۹۵۹ (۴)	Ionosphere
۰/۹۷۱ (۱)	۰/۹۶۷ (۳)	۰/۹۶۹ (۲)	۰/۹۴۹ (۴)	Pima Indians Diabetes
۰/۹۶۵ (۳)	۰/۹۶۷ (۲)	۰/۹۶۹ (۱)	۰/۹۴۹ (۴)	Spambase
۰/۹۵۹ (۲)	۰/۹۵۶ (۳)	۰/۹۶۲ (۱)	۰/۹۵۴ (۴)	Sonar

در این آزمایش، روند تغییرات صحت دسته‌بندها در چهار مرحله بررسی شده است. مرحله اول مجموعه داده‌ها بدون نویز است، در مراحل دیگر با افزایش نویز، به ترتیب ۲۰، ۱۰ و ۳۰ درصدی به مجموعه داده‌ها عملکرد دسته‌بندها مقایسه شده است. همان‌طور که در شکل ۸ مشخص است، صحت تمام دسته‌بندها، با افزایش نویز به‌میزان قابل توجهی کاهش یافته است.

جدول ۳: میانگین صحت دسته‌بندها در فرضیه داده‌های بدون نویز (کلاس ۱-
کلاس اصلی)

KH-SVDD	R-SVDD[25]	DW-SVDD[24]	SVDD[13]	مجموعه داده‌ها
۰/۹۷۲ (۱)	۰/۹۶۷ (۳)	۰/۹۶۹ (۲)	۰/۹۲۶ (۴)	Breast Cancer
۰/۹۸۱ (۲)	۰/۹۷۹ (۳)	۰/۹۸۲ (۱)	۰/۹۱۶ (۴)	Ionosphere
۰/۹۵۰ (۲)	۰/۹۴۹ (۳)	۰/۹۵۱ (۱)	۰/۸۹۵ (۴)	Pima Indians Diabetes
۰/۹۶۱ (۳)	۰/۹۶۲ (۱)	۰/۹۶۲ (۲)	۰/۸۷۹ (۴)	Spambase
۰/۹۴۲ (۱)	۰/۹۴۰ (۳)	۰/۹۴۲ (۲)	۰/۸۵۱ (۴)	Sonar

در جدول ۳ رتبه هر دسته‌بند داخل پرانتز نشان داده شده است. دسته‌بند پیشنهادی KH-SVDD، DW-SVDD و R-SVDD به‌علت استفاده از پارامتر تراکم در دسته‌بندی، نتایج بهتری را نسبت به دسته‌بند SVDD به‌دست آورده‌اند. هم‌چنین صحت دسته‌بندهای فوق نزدیک هم می‌باشند. جدول ۴ میانگین نرخ مثبت صحیح دسته‌بندها را در فرضیه اول نشان می‌دهد.

جدول ۴: میانگین نرخ مثبت صحیح دسته‌بندها در فرضیه داده‌های بدون نویز (کلاس ۱- کلاس اصلی)

KH-SVDD	R-SVDD[25]	DW-SVDD[24]	SVDD[13]	مجموعه داده‌ها
۰/۹۷۲ (۱)	۰/۹۶۷ (۳)	۰/۹۶۹ (۲)	۰/۹۴۱ (۴)	Breast Cancer
۰/۹۸۸ (۱)	۰/۹۸۲ (۳)	۰/۸۸۵ (۲)	۰/۹۲۹ (۴)	Ionosphere
۰/۹۵۸ (۲)	۰/۹۵۴ (۳)	۰/۹۶۳ (۱)	۰/۹۰۷ (۴)	Pima Indians Diabetes
۰/۹۶۳ (۳)	۰/۹۶۷ (۲)	۰/۹۶۹ (۱)	۰/۸۹۱ (۴)	Spambase
۰/۹۵۷ (۲)	۰/۹۵۹ (۱)	۰/۹۵۱ (۳)	۰/۸۶۲ (۴)	Sonar

جدول ۵ میانگین صحت دسته‌بندها و جدول ۶ میانگین نرخ مثبت صحیح در فرضیه اول (داده‌های بدون نویز) براساس آموزش دسته‌بندها با استفاده از کلاس ۲ به‌عنوان کلاس اصلی را نشان می‌دهند.

دسته‌بند پیشنهادی KH-SVDD در بیشتر آزمایش‌های داده‌های بدون نویز، از نظر صحت و نرخ مثبت صحیح رتبه اول و دوم را دارا است.

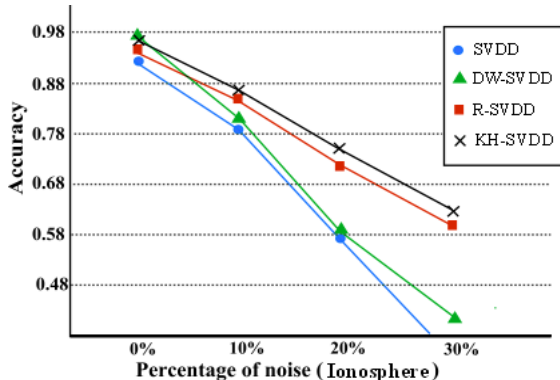
جدول ۷ نتایج صحت دسته‌بندها و جدول ۸ نتایج نرخ مثبت صحیح، در فرضیه دوم (داده‌های دارای ۱۰ درصد نویز) را نشان می‌دهند (کلاس ۱ به‌عنوان کلاس اصلی فرض شده است).

جدول ۹ نتایج صحت دسته‌بندها و جدول ۱۰ نتایج نرخ مثبت صحیح، در فرضیه دوم (داده‌های دارای ۱۰ درصد نویز) را نشان می‌دهند. در این آزمایش کلاس ۲ به‌عنوان کلاس اصلی فرض شده است.

جدول ۱۰: میانگین نرخ مثبت صحیح دسته‌بندها در فرضیه داده‌های دارای نویز (کلاس ۲- کلاس اصلی)

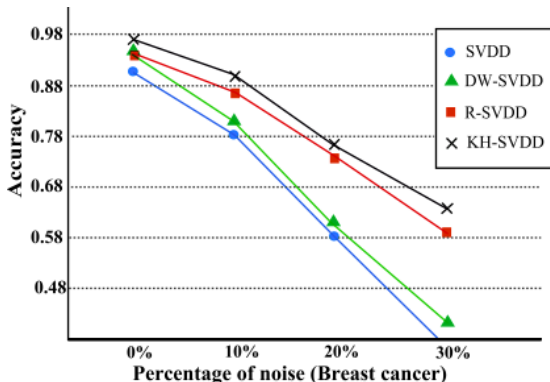
KH-SVDD	R-SVDD[25]	DW-SVDD[24]	SVDD[13]	مجموعه داده‌ها
۰/۹۱ (۱)	۰/۹۰ (۲)	۰/۸۸ (۳)	۰/۸۶ (۴)	Breast Cancer
۰/۸۷ (۲)	۰/۹۰ (۱)	۰/۸۵ (۳)	۰/۸۱ (۴)	Ionosphere
۰/۹۱ (۱)	۰/۸۹ (۲)	۰/۸۸ (۳)	۰/۸۱ (۴)	Pima Indians Diabetes
۰/۸۹ (۱)	۰/۸۸ (۲)	۰/۸۴ (۳)	۰/۷۷ (۴)	Spambase
۰/۹۲ (۱)	۰/۹۰ (۲)	۰/۸۷ (۳)	۰/۸۳ (۴)	Sonar

در شکل ۱۰، تغییرات صحت دسته‌بندها با افزایش نویز بر روی مجموعه داده‌های Pima indians diabetes نشان داده شده است. با افزایش ۳۰ درصدی نویز در مجموعه داده‌های فوق، عملکرد دسته‌بندهای KH-SVDD ۳۷ درصد، R-SVDD ۴۱ درصد، DW-SVDD ۵۹ درصد و SVDD ۶۱ درصد تقریباً کاهش یافته است.



شکل ۸: تغییرات صحت دسته‌بندها با افزایش نویز (Ionosphere)

در شکل ۱۱ برای مجموعه داده‌های Spambase، کاهش تقریبی عملکرد KH-SVDD ۳۸ درصد، R-SVDD ۴۰ درصد، DW-SVDD ۵۶ درصد و SVDD ۵۸ درصد با افزایش ۳۰ درصدی نویز، به دست آمده است.



شکل ۹: تغییرات صحت دسته‌بندها با افزایش نویز (Breast)

در شکل ۱۲، تغییرات صحت دسته‌بندها با افزایش نویز بر روی مجموعه داده‌های Sonar نشان داده شده است. با افزایش ۳۰ درصدی نویز، صحت دسته‌بندهای KH-SVDD ۳۴ درصد، R-SVDD ۳۵

جدول ۷: میانگین صحت دسته‌بندها در فرضیه داده‌های دارای نویز (کلاس ۱- کلاس اصلی)

KH-SVDD	R-SVDD[25]	DW-SVDD[24]	SVDD[13]	مجموعه داده‌ها
۰/۸۴ (۲)	۰/۸۵ (۱)	۰/۸۱ (۳)	۰/۷۹ (۴)	Breast Cancer
۰/۸۷ (۱)	۰/۸۶ (۲)	۰/۸۱ (۳)	۰/۷۸ (۴)	Ionosphere
۰/۸۶ (۱)	۰/۸۴ (۲)	۰/۸۲ (۳)	۰/۷۵ (۴)	Pima Indians Diabetes
۰/۸۷ (۱)	۰/۸۵ (۲)	۰/۸۳ (۳)	۰/۷۲ (۴)	Spambase
۰/۸۵ (۱)	۰/۸۳ (۲)	۰/۸۱ (۳)	۰/۷۳ (۴)	Sonar

در شکل ۸ تغییرات صحت دسته‌بندها با افزایش نویز بر روی مجموعه داده‌های Ionosphere نشان داده شده است. شیب خط عملکرد دسته‌بند پیشنهادی (KH-SVDD) و R-SVDD ملایم‌تر از دو روش دیگر است. برای مجموعه داده‌های Ionosphere، اختلاف درصد صحت تشخیص نمونه‌های اصلی و پرت در دسته‌بندی پیشنهادی بین دو بازه نویز صفر و ۳۰ درصد، تقریباً ۳۲ درصد است. این مقدار برای دسته‌بند R-SVDD، تقریباً ۳۷ درصد است.

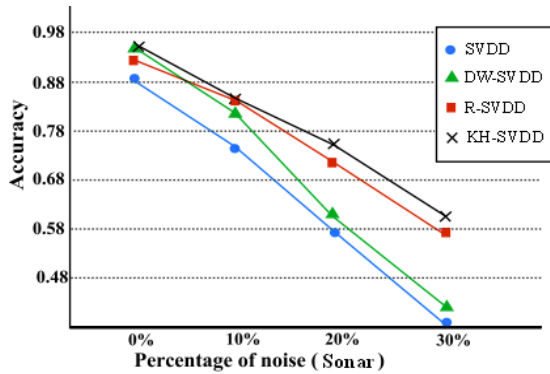
جدول ۸: میانگین نرخ مثبت صحیح دسته‌بندها در فرضیه داده‌های دارای نویز (کلاس ۱- کلاس اصلی)

KH-SVDD	R-SVDD[25]	DW-SVDD[24]	SVDD[13]	مجموعه داده‌ها
۰/۸۷ (۲)	۰/۸۹ (۱)	۰/۸۴ (۳)	۰/۸۲ (۴)	Breast Cancer
۰/۸۹ (۲)	۰/۹۰ (۱)	۰/۸۴ (۳)	۰/۸۰ (۴)	Ionosphere
۰/۸۸ (۱)	۰/۸۵ (۲)	۰/۸۲ (۳)	۰/۷۸ (۴)	Pima Indians Diabetes
۰/۹۰ (۱)	۰/۸۷ (۲)	۰/۸۵ (۳)	۰/۷۷ (۴)	Spambase
۰/۸۸ (۱)	۰/۸۶ (۲)	۰/۸۴ (۳)	۰/۷۶ (۴)	Sonar

در شکل ۹، تغییرات صحت دسته‌بندها با افزایش نویز بر روی مجموعه داده‌های Breast cancer نشان داده شده است. افزایش ۳۰ درصدی نویز در مجموعه داده‌های فوق، باعث کاهش تقریبی عملکرد دسته‌بندها به میزان ۳۰ درصد در KH-SVDD، ۳۴ درصد در R-SVDD، ۵۲ درصد در DW-SVDD و کاهش ۵۵ درصدی در دسته‌بند SVDD شده است.

جدول ۹: میانگین صحت دسته‌بندها در فرضیه داده‌های دارای نویز (کلاس ۲- کلاس اصلی)

KH-SVDD	R-SVDD[25]	DW-SVDD[24]	SVDD[13]	مجموعه داده‌ها
۰/۸۷ (۱)	۰/۸۶ (۲)	۰/۸۴ (۳)	۰/۸۰ (۴)	Breast Cancer
۰/۸۶ (۲)	۰/۸۷ (۱)	۰/۸۵ (۳)	۰/۷۹ (۴)	Ionosphere
۰/۸۸ (۱)	۰/۸۷ (۲)	۰/۸۶ (۳)	۰/۷۸ (۴)	Pima Indians Diabetes
۰/۸۸ (۱)	۰/۸۷ (۲)	۰/۸۵ (۳)	۰/۷۹ (۴)	Spambase
۰/۸۹ (۱)	۰/۸۶ (۲)	۰/۸۵ (۳)	۰/۷۹ (۴)	Sonar



شکل ۱۲: تغییرات صحت دسته‌بندی با افزایش نویز (Sonar)

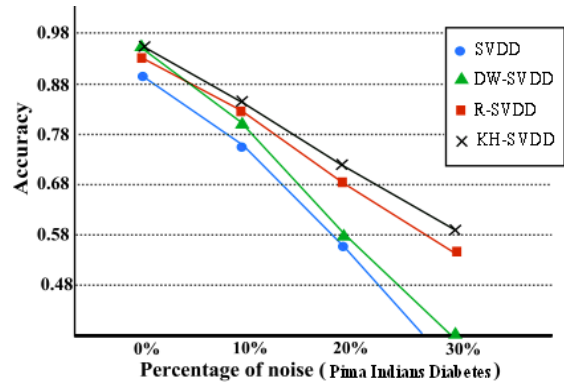
۶- نتیجه‌گیری و کارهای آینده

در این مقاله، یک دسته‌بند تک کلاسه برای داده‌های نویزی ارائه شد. موقعیت مرکز آبرگره در فضای ویژگی‌های داده‌ها با استفاده از الگوریتم SVDD و الگوریتم گروه میگوی آشوبی محاسبه شد. هم‌چنین برای کاهش تأثیر داده‌های نویزی بر مرکز آبرگره، از تراکم محلی نقاط و فاصله آنها تا مرکز آبرگره به‌عنوان یک تابع جریمه استفاده شده است. در دسته‌بند KH-SVDD، سعی بر تشخیص داده‌های نویزی در مجموعه داده‌های اصلی شده است. فرض کلی که در این روش وجود دارد این است که تنها مجموعه داده‌های اصلی موجود است و دانشی درباره داده‌های پرت وجود ندارد. نتایج آزمایش‌ها نشان می‌دهد روش پیشنهادی نسبت به روش‌های مشابه دیگر، از صحت بالاتری در تشخیص داده‌های پرت در مجموعه داده‌های نویزی، برخوردار است. کارهای آینده برای تشخیص داده‌های نویزی، می‌توان علاوه بر تراکم از پارامترهای آماری مانند میانگین، مد و انحراف معیار نیز استفاده کرد.

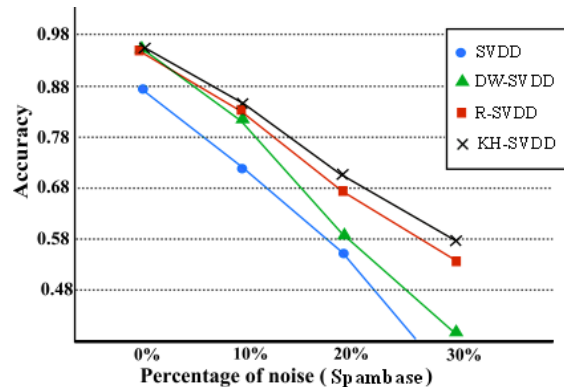
مراجع

- [۱] محمدعلی زارع چاهوکی و سیدحمیدرضا محمدی، «بهینه‌سازی هسته‌های چندگانه در ماشین بردار پشتیبان جفتی برای کاهش شکاف معنایی تشخیص صفحات فریب‌آمیز»، مجله مهندسی برق دانشگاه تبریز، شماره ۴ جلد ۴۶، ۱۳۵-۱۴۵، ۱۳۹۵.
- [۲] محمدمامیر عباسیان و حسین نظام‌آبادی‌پور، «الگوریتم جستجوی گرانشی چندهدفه مبتنی بر مرتب‌سازی جبهه‌های مغلوب‌نشده»، مجله مهندسی برق دانشگاه تبریز، شماره ۱ جلد ۴۱، ۶۸-۸۰، ۱۳۹۰.
- [۳] سیدحسین غفاریان، هادی صدوقی یزدی و یونس الله‌یاری، «دسته‌بند تک‌کلاسه گرانش‌گرای مبتنی بر ماشین بردار پشتیبان»، نشریه مهندسی برق و مهندسی کامپیوتر ایران، سال ۱۰، شماره ۲، ۱۳۹۱.
- [۴] وحیده منعمی‌زاده و جواد حمیدزاده، «جستجوی k نزدیک‌ترین همسایه تقریبی به روش ترکیب خطی»، نشریه مهندسی برق و مهندسی کامپیوتر ایران، آماده انتشار.

درصد، ۵۲ DW-SVDD و ۵۲ درصد به‌طور تقریبی کاهش یافته است.



شکل ۱۰: تغییرات صحت دسته‌بندی با افزایش نویز (Diabetes)



شکل ۱۱: تغییرات صحت دسته‌بندی با افزایش نویز (Spambase)

نتایج آزمایش‌ها نشان می‌دهد که دسته‌بند پیشنهادی نسبت به دسته‌بندهای بررسی‌شده دیگر، در هنگام افزایش نویز نتایج بهتری را به‌دست آورده است. کارآیی دو دسته‌بند SVDD و DW-SVDD به‌علت نداشتن رویکردی برای مقابله با داده‌های نویزی، در زمان حضور نویز دچار افت شدیدی شده‌اند. روش R-SVDD و KH-SVDD با توجه به نداشتن رویکردی برای مقابله با نویز، نسبت به دو روش دیگر کارآیی بهتری را در آزمایش‌ها کسب کردند.

یکی از مشکلات روش پیشنهادی افزایش زمان اجرای الگوریتم در مرحله آموزش دسته‌بند است. یافتن مرکز مناسب برای کلاس اصلی داده‌ها توسط الگوریتم میگو انجام می‌شود که یک سربار اضافی بر روی دسته‌بند ایجاد می‌کند.

- [18] K. Lee, D. Kim, K. H. Lee and D. Lee, "Density-induced support vector data description," *Neural Networks, IEEE Transactions on*, vol. 18, pp. 284–289, 2007.
- [19] C. K. Wang, Y. Ting, Y. H. Liu and G. Hariyanto, "A Novel Approach to Generate Artificial Outliers for support Vector Data Description," *IEEE International Symposium on Industrial Electronics (ISIE)*, Korea, pp. 2202-2207, 2009.
- [20] H. W. Cho, "Data description and noise filtering based detection with its application and performance comparison," *Expert systems with Applications*, vol. 36, no. 1, pp. 434-441, 2009.
- [21] S. M. Guo, L. C. Chen and J. S. Tsai, "A boundary method for outlier detection based on support vector domain description," *Pattern Recognition*, vol. 42, pp. 77-83, 2009.
- [22] G. X. Huang, H. F. Chen and F. Yin, "Improved support vector data description," *International Conference on Machine Learning and Cybernetics*, vol. 3, pp. 1459-1463, 2010.
- [23] B. Liu, Y. Xiao, L. Cao, Z. Hao and F. Deng, "SVDD-based outlier detection on uncertain data," *Knowledge and Information Systems*, vol. 34, pp. 597-618, 2013.
- [24] M. Cha, J. Kim and J. Baek, "Density weighted support vector data description," *Expert Systems with Applications*, vol. 41, pp. 3343–3350, 2014.
- [25] G. Chen, X. Zhang, Z. Wang and F. Lia, "Robust support vector data description for outlier detection with noise or uncertain data," *Knowledge-Based Systems*, vol. 90, pp. 129–137, 2015.
- [26] S. Kim, Y. Choi and M. Lee, "Deep learning with support vector data description," *Neurocomputing*, vol. 165, pp. 111–117, 2015.
- [27] G. Wang, L. Guo, A. Gandomi, G. Hao and H. Wang, "Chaotic Krill Herd algorithm," *Information Sciences*, vol. 274, pp. 17–34, 2014.
- [28] A. Asuncion and D. Newman, *UCI Machine Learning Repository*. University of California, School of Information and Computer Science, Irvine, CA, 2013.
- [29] S.S. Khan, J. Hoey and D. Lizotte, "Bayesian multiple imputation approaches for one-class classification," *Advances in Artificial Intelligence*, pp. 331–336, 2012.
- [30] B. Liu, Y. Xiao and Z. Hao, "An efficient approach for outlier detection with imperfect data labels," *IEEE Trans Knowl. Data Eng.*, vol. 26, pp. 1602-1616, 2014.
- [5] S. S. Khan and M. G. Madden, "A survey of recent trends in one class classification." *Artificial Intelligence and Cognitive Science*, vol. 6206, pp. 188-197, 2010.
- [6] A. Wenjuan, M. Liang and H. Liu, "An improved one-class support vector machine classifier for outlier detection," *Mechanical Engineering Science*, vol. 229, pp. 580-588, 2015.
- [7] S. S. Khan and M. G. Madden, "One-Class Classification: Taxonomy of Study and Review of Techniques," *The Knowledge Engineering Review*, vol. 29, pp. 1-30, 2014.
- [8] S. Kang, S. Cho and P. Kang, "Multi-class classification via heterogeneous ensemble of one-class classifiers," *Engineering Applications of Artificial Intelligence*, vol. 43, pp. 35–43, 2015.
- [9] L. Zhang, L. Xingning, W. Bangjun and H. Shuping, "Similarity learning based on multiple support vector data description," *Neural Networks (IJCNN)*, pp. 1-7, 2015.
- [10] D. M. Tax and R.P. Duin, "Uniform object generation for optimizing one-class classifiers," *The Journal of Machine Learning Research*, vol. 2, pp. 155-173, 2002.
- [11] R. Sadeghi and J. Hamidzadeh, "Automatic Support Vector Data Description," *Soft Computing*, 2016, DOI: 10.1007/s00500-016-2317-5.
- [12] V. H. Moghaddam, and J. Hamidzadeh, "New Hermite orthogonal polynomial kernel and combined kernels in Support Vector Machine classifier," *Pattern Recognition*, vol. 60, pp. 921-935, 2016.
- [13] D. M. Tax and R.P. Duin, "Support vector data description," *Machine Learning*, vol. 54, pp. 45–66, 2004.
- [14] J. Bootkrajang, "A generalised label noise model for classification in the presence of annotation errors," *Neurocomputing*, vol. 192, pp. 61–71, 2016.
- [15] J. Hamidzadeh, R. Monsefi and H. SadoghiYazdi, "IRAHC: Instance Reduction Algorithm using Hyperrectangle Clustering," *Pattern Recognition*, vol. 48, pp.1878-1889, 2015.
- [16] J. Hamidzadeh, R. Monsefi and H. SadoghiYazdi, "LMIRA: Large Margin Instance Reduction Algorithm," *Neurocomputing*, vol. 145, pp. 477-487, 2014.
- [17] S. Y. Xia, Z. Xiong, Y. He, K. Li, L. M. Dong and M. Zhang, "Relative density-based classification noise detection," *Optik International Journal for Light and Electron Optics*, vol. 125, pp. 6829–6834, 2014.

زیر نویس ها

¹ Outlier detection

² One-class classification

³ Classification

⁴ Krill herd support vector data description

⁵ Support vector data description

⁶ Kernel trick

⁷ Parzen

⁸ Post-processing

⁹ Likelihood

¹⁰ Overfitting

¹¹ 10-fold cross validation

¹² Confusion