



به نام خدا



دانشکده مهندسی
گروه کامپیوتر
آزمایشگاه شناسایی الگو

پایان نامه کارشناسی ارشد

تلفیق توصیفگرهای داده تک کلاسی

مبتنی بر بردار پشتیبان

نگارش

نشاط صالحی

استاد راهنما

دکتر هادی صدوقی یزدی

استاد مشاور

دکتر عباس قائمی بافقی

به نام خدا



دانشکده مهندسی گروه کامپیوتر
آزمایشگاه شناسایی الگو
پایان نامه کارشناسی ارشد

عنوان:

تلفیق توصیفگرهای داده تک کلاسی
مبتنی بر بردار پشتیبان

دانشجو:

نشاط صالحی

کمیته ممتحنین:

استاد راهنما: دکتر هادی صدوقی یزدی امضا:

استاد مشاور: دکتر عباس قائمی بافقی امضا:

استاد داور: دکتر محمدباقر نقیبی سیستانی امضا:

استاد داور و نماینده کمیته تحصیلات تکمیلی: دکتر رضا منصفی امضا:

تعهدنامه

اینجانب نشاط صالحی دانشجوی دوره کارشناسی ارشد رشته مهندسی کامپیوتر هوش مصنوعی دانشکده مهندسی دانشگاه فردوسی مشهد نویسنده پایان نامه تلفیق توصیفگرهای داده تک کلاسی مبتنی بر بردار پشتیبان تحت راهنمایی استاد ارجمند جناب آقای دکتر هادی صدوقی یزدی و جناب آقای دکتر عباس قائمی بافقی متعهد می شوم:

- تحقیقات در این پایان نامه توسط اینجانب انجام شده و از صحت و اصالت برخوردار است.
- در استفاده از نتایج پژوهش‌های محققان دیگر به مرجع مورد استفاده استناد شده است.
- مطالب مندرج در پایان نامه تاکنون توسط خود و یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است.
- کلیه حقوق معنوی این اثر متعلق به دانشگاه فردوسی مشهد می‌باشد و مقالات مستخرج با نام «دانشگاه فردوسی مشهد» و یا "Ferdowsi University of Mashhad" به چاپ خواهد رسید.
- حقوق معنوی تمام افرادی که در به دست آمدن نتایج اصلی پایان نامه تأثیرگذار بوده‌اند در مقالات مستخرج از رساله رعایت شده است.
- در کلیه مراحل انجام این پایان نامه، در مواردی که از موجود زنده (یا بافت‌های آن‌ها) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است.
- در کلیه مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است، اصل رازداری، ضوابط و اصول اخلاق انسانی رعایت شده است.

تاریخ

امضای دانشجو

نشاط صالحی

مالکیت نتایج و حق نشر

- کلیه حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه‌های رایانه‌ای، نرم‌افزارها و تجهیزات ساخته شده) متعلق به دانشگاه فردوسی مشهد می‌باشد. این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات و نتایج موجود در پایان نامه بدون ذکر مرجع مجاز نمی‌باشد.



تقدیم بہ

پدرم بہ استواری کوہ

روح پاک مادرم بہ زلالی چشمہ



ہمسرم بہ صمیمیت باران

و اساتید کرامت قدم...

تقدیر و تشکر

سپاس و آفرین ایزد جهان آفرین راست. سپاس خدای را که حق ستایش برتر از حد ستایشگران است و نعمت‌هایش فوق اندیشه شمارشگران. سپاس خدای را که پای اندیشه تیزگام در راه شناسایی او لنگ است و سر فکر ژرف رو به دریای معرفتش بر سنگ.

از این رو بر خود لازم می‌دانم که سپاس گویم تمام کسانی را که از آن‌ها آموختم. تشکر می‌نمایم از استاد گرانقدرم جناب آقای دکتر هادی صدوقی یزدی که فروغ اندیشه‌شان زدایندهی ظلمت جهل در حرکت به سوی علم و آگاهی است و نیز ارج می‌نهم راهنمایی‌های استاد بزرگوار جناب آقای دکتر قائمی بافقی، که بزرگوارانه مرا راهنمایی نمودند.

سپاس خود را نثار همه معلمان و اساتیدم می‌دارم که از نخستین مراحل تحصیل تاکنون، مشوق، معلم و راهنمای من در مسیر علم و زندگی بودند.

سپاس می‌گویم خانواده عزیزم را که همیشه پشتیبان و مشوقم بوده‌اند و شرایطی را فراهم آوردند تا در محیطی مطلوب و پر از صفا و آرامش مراتب تحصیل را طی نمایم. در انتها قدردانی می‌کنم زحمات همسر و جملگی دوستانی را که در تهیه این اثر سهمی داشته‌اند، امید که برای پویندگان راه دانش مثمر ثمر واقع گردد.

چکیده

توصیف‌گر داده مبتنی بر بردار پشتیبان (SVDD)، یک طبقه‌بند باناظر تک کلاسه است. هدف این طبقه‌بند مرزی، بهینه کردن حجم دایره (ابکره) اطراف مجموعه هدف خطی یا غیرخطی می‌باشد. حداقل پیچیدگی زمانی این گونه طبقه‌بندها، $O(n^3)$ است؛ در نتیجه با افزایش تعداد نمونه‌ها، مسئله برای مجموعه داده‌های حجیم کارایی خود را از دست می‌دهد. هدف اصلی این پایان‌نامه، توسعه SVDD، به منظور ایجاد امکان استفاده از آن در کاربردهای حجیم و افزایش سرعت بدون کاهش کارایی یادگیری در توصیف داده‌ها است. این هدف، با بهره‌گیری از روش‌های تلفیق طبقه‌بندها به منظور افزایش سرعت یادگیری و کاهش نسبی میزان خطا در توصیف داده‌ها برآورده شده است؛ که نتیجه‌ی آن تلفیق طبقه‌بندها با وزن دهی مبتنی بر کارایی آن‌ها می‌باشد. در بند اول روش پیشنهادی، طبقه‌بند SVDD، ابتدا با روش AdaBoost ترکیب گردیده و سپس در بند دوم، ضرایب لاگرانژ طبقه‌بند SVDD، با استفاده از روش تخمین مبتنی بر بیشینه کردن امیدریاضی (EM)، محاسبه شده است که نتیجه‌ی آن دو نوع طبقه‌بند ترکیبی همگرا است. دو روش پیشنهادی از نظر سرعت و دقت با سایر الگوریتم‌های تک کلاسه SVDD مبتنی بر کرنل از جمله FSVDD و طبقه‌بند افزایشی (Inc-SVDD) و جعبه‌ابزار LibSVM-SVDD بر روی مجموعه داده‌های استاندارد UCI و مجموعه داده IDS آزمایش و مقایسه شده‌اند. نتایج حاصل، برتری روش‌های پیشنهادی را از نظر کارایی و سرعت در مرحله آموزش و تست نشان می‌دهد.

کلیدواژه: توصیف‌گر داده مبتنی بر بردار پشتیبان - داده‌های حجیم - ترکیب طبقه‌بندها - روش آدابوست - روش تخمین مبتنی بر بیشینه کردن امیدریاضی (EM).

فهرست مطالب

I.....	تعهدنامه
V.....	تقدیر و تشکر
VI.....	چکیده
XI.....	فهرست شکل‌ها
XIII.....	فهرست جداول
XIV.....	فهرست علائم اختیاری
XVI.....	فهرست اختصارات

I **فصل (۱) مقدمه**

۴.....	۱-۱- راهکارهای تسریع
۴.....	۱-۱-۱- سطح داده
۴.....	۱-۱-۲- سطح ویژگی
۵.....	۱-۱-۳- سطح طبقه‌بند
۵.....	۱-۱-۴- سطح ترکیب طبقه‌بندها
۶.....	۲-۱- تبیین موضوع پایان‌نامه
۸.....	۳-۱- شمای کلی پایان‌نامه

۱۰..... **فصل (۲) مفاهیم اولیه مورد نیاز**

۱۱.....	۱-۲- تابع کرنل
۱۱.....	۱-۱-۲- مزایای تابع کرنل
۱۲.....	۲-۲- معرفی طبقه‌بند SVDD و پایه ریاضی آن
۱۳.....	۱-۲-۲- فرمول‌بندی مسئله

۱۴	۲-۲-۲ فرمول‌بندی مسئله با استفاده از تابع کرنل
۱۵	۳-۲-۲ بررسی ویژگی‌های طبقه‌بند SVDD
۱۶	۳-۲-۳ معرفی الگوریتم تقویت انطباقی (آدابوست)
۱۸	۱-۳-۲ همگرایی آدابوست
۲۰	۲-۳-۲ ویژگی‌های آدابوست
۲۰	۴-۲-۳ معرفی الگوریتم انتخاب چرخ رولت
۲۲	۵-۲-۳ معرفی روش بیشینه کردن امیدریاضی
۲۴	۱-۵-۲ همگرایی روش بیشینه کردن امید ریاضی
۲۶	۲-۵-۲ ویژگی‌ها و معایب الگوریتم EM
۲۷	۶-۲-۳ معرفی مدل ترکیبی گوسی (GMM)
۲۸	۷-۲-۳ خلاصه

فصل ۳ (۳) مروری بر کارهای پیشین..... ۲۹

۳۰	۱-۳-۱ روش‌های تسریع در سطح داده
۳۱	۱-۱-۳ تسریع توسط استخراج نمونه‌های مرزی [LIA '09]
۳۴	۲-۱-۳ روش K دورترین - همسایه [Xia '10]
۳۷	۳-۱-۳ تسریع بر اساس گسسته‌سازی و ترکیب داده‌ها [Luo '10]
۳۹	۴-۱-۳ طبقه‌بند افزایشی SVDD [Hua '11]
۴۰	۲-۳-۲ روش‌های تسریع در سطح ویژگی
۴۲	۱-۲-۳ انتخاب ویژگی برای طبقه‌بند ماشین بردار پشتیبان (FS_SFS) [Liu '06]
۴۷	۲-۲-۳ سایر روش‌های کاهش ویژگی
۴۸	۳-۳-۳ روش‌های تسریع در سطح طبقه‌بندها
۴۹	۱-۳-۲ تسریع در زمان آموزش
۵۶	۲-۳-۲ تسریع در زمان آزمون
۶۱	۴-۳-۴ روش‌های تسریع در سطح ترکیب طبقه‌بندها
۶۳	۱-۴-۳ پیاده‌سازی SVM با شبکه عصبی مؤلفه‌ای [Hua '05]
۶۶	۲-۴-۲ پیاده‌سازی SVDD با الگوریتم ژنتیک [Tav '07]
۶۸	۳-۴-۲ سایر روش‌های سطح ترکیب طبقه‌بند

۶۸.....۳-۵-خلاصه

فصل ۴) روش پیشنهادی.....۷۱

۷۲.....۴-۱-روش اول: بهبود عملکرد SVDD با استفاده از روش آدابوست

۷۳.....۴-۱-۱- معرفی روش پیشنهادی ABSVDD

۷۵.....۴-۱-۲- تشریح الگوریتم مرحله آموزش

۷۶.....۴-۱-۳- تشریح الگوریتم مرحله آزمون

۷۷.....۴-۱-۴- آنالیز بیش برآزش روش پیشنهادی

۷۹.....۴-۱-۵- آنالیز مرتبه زمانی روش پیشنهادی ABSVDD

۸۰.....۴-۱-۶- آنالیز همگرایی روش پیشنهادی ABSVDD

۸۱.....۴-۱-۷- مثالی از روند آموزش روش پیشنهادی ABSVDD

۸۳.....۴-۱-۸- تفسیر ضرایب وزنی طبقه‌بندهای ABSVDD

۸۴.....۴-۲-روش دوم: بهبود عملکرد SVDD با استفاده روش EM

۸۴.....۴-۲-۱- معرفی روش پیشنهادی EMSVDD

۹۰.....۴-۲-۲- تشریح الگوریتم مرحله آموزش

۹۱.....۴-۲-۳- تشریح الگوریتم مرحله آزمون

۹۱.....۴-۲-۴- آنالیز مرتبه زمانی روش پیشنهادی EMSVDD

۹۳.....۴-۲-۵- آنالیز همگرایی روش پیشنهادی

۹۳.....۴-۲-۶- مثالی از روند آموزش روش پیشنهادی EMSVDD

۹۵.....۴-۲-۷- تفسیر ضرایب وزنی طبقه‌بندهای EMSVDD

۹۵.....۴-۳-خلاصه

فصل ۵) آزمایش و ارزیابی.....۹۷

۹۸.....۵-۱-بستر ارزیابی

۹۹.....۵-۲-معیارهای ارزیابی

۹۹.....۵-۱-۲- معیارهای حاصل از ماتریس در هم ریختگی دو کلاس

۱۰۱.....۵-۲-۲- سایر معیارها

۱۰۱.....۵-۳-مجموعه داده‌ها

۱۰۱.....	مجموعه داده‌های مصنوعی.....	۱-۳-۵
۱۰۲.....	مجموعه داده‌های معمولی و بزرگ.....	۲-۳-۵
۱۰۴.....	مجموعه داده‌ی تشخیص نفوذ.....	۳-۳-۵
۱۰۵.....	روش‌های رقیب.....	۴-۳-۵
۱۰۶.....	مقایسه طبقه‌بندهای پیشنهادی ABSVDD و EMSVDD با روش‌های رقیب.....	۴-۵
۱۰۶.....	بررسی نتایج روش‌های پیشنهادی بر روی مجموعه داده‌های مصنوعی.....	۲-۴-۵
۱۰۹.....	آزمایش بر مجموعه داده‌های کوچک و متوسط.....	۳-۴-۵
۱۱۳.....	آزمایش بر مجموعه داده‌های بزرگ.....	۴-۴-۵
۱۱۶.....	بررسی نتایج روش‌های پیشنهادی در تشخیص نفوذ به شبکه‌های کامپیوتری.....	۵-۴-۵
۱۱۸.....	بررسی میزان تغییرات دقت نسبت به پارامتر C.....	۶-۴-۵
۱۱۹.....	خلاصه.....	۵-۵

فصل ۶) نتیجه‌گیری و توصیه‌های آتی..... ۱۲۱

۱۲۲.....	مروری بر روش‌های پیشنهادی.....	۱-۶
۱۲۳.....	توصیه‌های آتی.....	۲-۶
۱۲۳.....	اعمال روش بر روی داده‌های جریان‌ی.....	۱-۲-۶
۱۲۳.....	تخمین پارامترهای تابع کرنل گوسی در روش پیشنهادی EMSVDD.....	۲-۲-۶
۱۲۳.....	سایر کارهای ممکن.....	۳-۲-۶
۱۲۴.....	خلاصه.....	۳-۶

جمع بندی..... ۱۲۵

۱۳۰.....	منابع و مأخذ.....	
----------	-------------------	--

فهرست شکل‌ها

- شکل ۱-۱: طبقه‌بند تک کلاسی مبتنی بر بردار پشتیبان. ۶
- شکل ۲-۲: داده‌های جداناپذیر خطی توسط تابع کرنل تبدیل به داده‌های جداپذیر خطی شده‌اند [DTR] ۱۱
- شکل ۳-۲: داده‌های هدف توسط یک کره با شعاع R و مرکز a محدود شده‌اند. ۱۲
- شکل ۴-۲: چرخ رولت برای پنج داده با توجه به میزان برازندگی آنها، مشخص شده است [Mat] ۲۱
- شکل ۵-۲: فلوجارت روش بیشینه کردن امیدریاضی [Kha '11] ۲۴
- شکل ۱-۳: توزیع نمونه‌های انتخاب شده با پارامتر ϵ ‌های مختلف [LIA '09] ۳۳
- شکل ۲-۳: توزیع نمونه‌های انتخاب شده با پارامتر r و ϵ ‌های مختلف [LIA '09] ۳۴
- شکل ۳-۳: مقایسه دقت برای ϵ ‌های مختلف و تحت $r=0.01$ و $r=0.04$ [LIA '09] ۳۴
- شکل ۴-۳: ترکیب فیلترینگ معمولی و فیلترینگ حلقوی [Liu '06] ۴۴
- شکل ۵-۳: مقایسه سه روش Osuna، Chuncking و SMO [Pla '99] ۵۰
- شکل ۶-۳: خصوصیات هندسی SVDD [Liu '10] ۵۷
- شکل ۷-۳: ساختار TLFN با L ماژول کوانتایزر و L ماشین آموزش [Hua '05] ۶۳
- شکل ۸-۳: کدینگ و نمایش کروموزوم‌های مسئله SVDD [Tav '07] ۶۷
- شکل ۱-۴: شکستن مجموعه داده‌ی اصلی به دو بخش مجزا و آموزش آن‌ها. ۷۸
- شکل ۲-۴: نمایش مراحل اجرای روش پیشنهادی ABSVDD ۸۲
- شکل ۳-۴: مقایسه میزان خطا و وزن هر طبقه‌بند پایه در روش ABSVDD ۸۳
- شکل ۴-۴: نمایش مراحل اجرای روش پیشنهادی EMSVDD ۹۴
- شکل ۱-۵: مجموعه داده‌های مصنوعی. ۱۰۲
- شکل ۲-۵: نمودار میله‌ای معیار دقت الگوریتم‌ها بر مجموعه داده‌های مصنوعی ۱۰۷
- شکل ۳-۵: نمودار دقت الگوریتم‌ها بر مجموعه داده‌های مصنوعی ۱۰۸
- شکل ۴-۵: نمودار میله‌ای زمان آموزش بر مجموعه داده‌های مصنوعی ۱۰۸
- شکل ۵-۵: نمودار دقت الگوریتم‌ها روی مجموعه داده‌های کوچک و متوسط ۱۱۱
- شکل ۶-۵: نمودار نرخ خطای مرحله یادگیری الگوریتم‌ها روی مجموعه داده‌های کوچک و متوسط ۱۱۱
- شکل ۷-۵: نمودار درصد بردارهای پشتیبان الگوریتم‌ها روی مجموعه داده‌های کوچک و متوسط ۱۱۲

- شکل ۵-۸: نمودار میله‌ای تعداد نمونه‌های بردار پشتیبان روی مجموعه داده‌های بزرگ ۱۱۴
- شکل ۵-۹: نمودار میله‌ای درصد قابلیت اعتماد به خروجی روی مجموعه داده‌های بزرگ ۱۱۴
- شکل ۵-۱۰: نمودار میله‌ای زمان آموزش روی مجموعه داده‌های بزرگ ۱۱۵
- شکل ۵-۱۱: نمودار خطای مجموعه داده‌های بزرگ طی تکرارهای متوالی روش پیشنهادی EMSVDD ۱۱۶
- شکل ۵-۱۲: نمودار میزان تغییرات دقت نسبت به پارامتر C ۱۱۹

فهرست جداول

- جدول ۱-۲: متداول ترین توابع کرنل ۱۲
- جدول ۱-۳: روش های موجود در سطح طبقه بند [Men '07] ۶۱
- جدول ۲-۳: مروری بر روش های تسریع طبقه بندهای مبتنی بر ماشین بردار پشتیبان ۷۰
- جدول ۱-۵: ماتریس در هم ریختگی ۹۹
- جدول ۲-۵: تعداد نمونه های مجموعه داده های مصنوعی در اندازه بزرگ ۱۰۱
- جدول ۳-۵: مجموعه داده های حقیقی ۱۰۳
- جدول ۴-۵: تعداد نمونه های مجموعه داده ی تشخیص نفوذ ۱۰۵
- جدول ۵-۵: نتایج مربوط به مجموعه داده های مصنوعی الگوریتم های تسریع ۱۰۷
- جدول ۶-۵: نتایج مربوط به مجموعه داده های با اندازه کوچک و معمولی، توسط الگوریتم های تسریع ۱۱۰
- جدول ۷-۵: نتایج مربوط به مجموعه داده های بزرگ توسط الگوریتم های تسریع ۱۱۳
- جدول ۸-۵: ماتریس در هم ریختگی روش ABSVDD برای مجموعه داده NSL-KDD ۱۱۶
- جدول ۹-۵: نتایج مربوط به معیار دقت مجموعه داده NSL-KDD ۱۱۷
- جدول ۱۰-۵: نتایج ارزیابی مجموعه داده NSL-KDD نسبت به معیارهای FAR, DR, CA ۱۱۸

فهرست علائم اختیاری

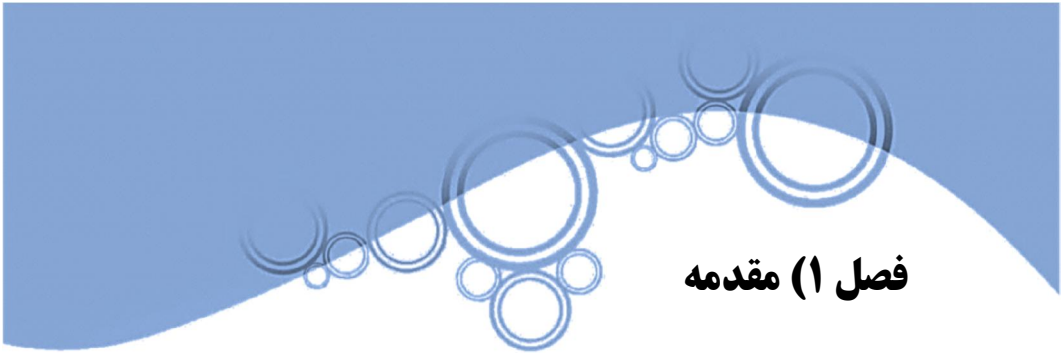
در لیست زیر به بیان متغیرهای به‌کاررفته در روش پیشنهادی در پایان‌نامه اشاره شده است.

عنوان	توضیح
R	شعاع ابرکره مرز توصیفگر داده
a	مرکز ابرکره مرز توصیفگر داده
ξ_i	متغیر لغزش به ازای هر نمونه
C	پارامتر تنظیم نسبت حجم ابرکره و میزان خطا
β_i	ضریب لاگرانژ به ازای هر نمونه
γ_i	ضریب لاگرانژ به ازای هر نمونه
$\varphi : x \rightarrow \varphi(x)$	تابع انتقال به فضای ویژگی (فضای با ابعاد بالاتر)
$K(x_i, x_j)$	تابع کرنل
σ	پارامتر تابع کرنل RBF
z	نمونه مورد آزمون
ω_i	وزن طبقه‌بند یا وزن ترکیب
h	طبقه‌بند
Z_t	فاکتور نرمال‌سازی
f_k	مقدار برازندگی نمونه k ام
P_k	احتمال انتخاب
$L(\theta)$	تابع درست‌نمایی
J	تابع حد پایین درست‌نمایی
θ	بردار پارامترهای نامعلوم در تابع چگالی احتمال
$M(\theta)$	تابع نگاشت از فضای پارامتر به خودش
$Jh(\theta)$	ماتریس ژاکوبی $d \times d$
$\tau_c(\theta; y)$	امیدریاضی تابع مشتق دوم لگاریتم درست‌نمایی نسبت به مشاهدات
$N(\mu_i, \Sigma_i)$	تابع توزیع نرمال
M	تعداد ترکیب از چند تابع توزیع نرمال یا چند طبقه‌بند

عنوان	توضیح
μ_i	بردار میانگین تابع توزیع نرمال
Σ_i	ماتریس کوواریانس تابع توزیع نرمال
d	بعد فضای ورودی‌ها
$p(X; \theta)$	تابع توزیع نمونه‌های X با پارامترهای θ
$SVDD(x; \theta)$	تابع توزیع SVDD برای نمونه‌های X با پارامترهای θ
$ SV $	تعداد بردارهای پشتیبان
k	درصدی از تعداد داده‌ها
$f_k(x)$	تابع تصمیم طبقه‌بند
$\delta(C_k x)$	تابع شبه توزیع احتمالی پسین
$F(x_n; \theta)$	تابع ترکیب چند گوسی
$\Lambda(X; \theta)$	تابع درست‌نمایی
$\lambda(X; \theta)$	تابع لگاریتم درست‌نمایی
$Q(\theta, \theta^s)$	تابع امیدریاضی شرطی تابع درست‌نمایی
NC	تعداد طبقه‌بندها

فهرست اختصارات

اختصار	توضیح
AdaBoost	Adaptive Boosting
EM	Expectation Maximization
FS_SFS	Filtered and Supported Sequential Forward Search
GMM	Gaussian Mixture Model
IDS	Intrusion Detection Systems
KFN-CBD	K-Farthest-Neighbors-Based Concept Boundary Determination
KKT	Karush–Kuhn–Tucker
LIBSVM	Library for Support Vector Machines
NSV	Non Support Vector
QP	Quadratic Programming
RBF	Radial Basis Function
SGD	Stochastic Gradient Descend
SLFNs	Single-hidden Layer Feed forward Networks
SMO	Sequential Minimal Optimization
SVNUB	Support Vector Not Upper Bound
SVDD	Support Vector Data Description
SVM	Support Vector Machine
SVUB	Support Vector Upper Bound
TLFN	Two-hidden Layer Feed-forward Network
UBSVs	Unbounded Support Vectors



فصل (۱) مقدمه

این فصل به معرفی مختصری از مسائل مطرح شده در پایان نامه می پردازد و شامل موارد زیر

می باشد:

✓ دلایل انتخاب موضوع پایان نامه و بیان اهمیت آن

✓ تبیین موضوع پایان نامه

✓ مروری کلی بر مطالب پایان نامه

مطالعه این فصل در جهت دادن به سایر مطالب پایان نامه مهم بوده و می تواند مفید واقع شود.

افزایش روزافزون حجم اطلاعات ذخیره شده و عدم قابلیت استفاده از آرشیوهای عظیم اطلاعات خام در تصمیم‌گیری‌ها، نیاز به فرآیندی که بتواند دانش موجود در این اطلاعات را کشف نماید هر روز نمایان‌تر می‌سازد. شناسایی الگو زمینه تحقیقاتی است که به طراحی سیستم‌هایی برای تشخیص الگوها در داده‌ها و توصیف داده‌ها می‌پردازد و ما را در رسیدن به درکی از ساختار داده‌ها و به دست آوردن دانش و تصمیم‌گیری یاری می‌نماید. توسعه آن به ۱۹۶۰ برمی‌گردد و با ارتباط با رشته‌های دیگر مانند آمار، مهندسی، هوش مصنوعی، علوم کامپیوتر، روانشناسی و... به عنوان زمینه‌ای کاملاً بین رشته‌ای [Web '99] شناخته می‌شود. شناسایی الگو کاربردهای بسیار متنوعی در علوم مختلف مانند بینایی ماشین، سنجش از راه دور، تشخیص‌های پزشکی، داده‌کاوی (در کاربردهایی مانند تحلیل ریسک، ارزیابی اعتبار، تحلیل فروش و...) و موارد بسیار متعدد دیگر دارد.

یکی از روش‌های شناسایی الگو که در حال حاضر به طور گسترده برای مسئله دسته‌بندی مورد استفاده قرار می‌گیرد، طبقه‌بندهای مبتنی بر بردار پشتیبان^۱ است [LEE '07a]. در سال‌های اخیر طبقه‌بندهای مبتنی بر بردار پشتیبان برای شناسایی و دسته‌بندی الگو و در کاربردهایی مانند پردازش تصاویر و ویدئو [Kha '12]، داده‌کاوی، کاربردهای پزشکی [Ji '08] و سیستم‌های تشخیص نفوذ^۲ به طور گسترده‌ای مورد استفاده قرار گرفته است [Gha '10]. نتایج حاصل از به‌کارگیری این روش‌ها در کاربردهای متفاوت، حاکی از کارایی بالای آن‌ها است.

این طبقه‌بندها دارای قابلیت‌های ارزشمندی هستند که آن‌ها را نسبت به دیگر روش‌های موجود برتر ساخته است [Tax '99]؛

(۱) آموزش خود مشکل بیشینه و کمینه محلی را ندارند.

(۲) دسته‌بندی‌کننده را با حداکثر تعمیم بنا می‌کنند.

(۳) ساختار و توپولوژی خود را به صورت بهینه تعیین می‌نمایند

¹ Support Vector Machines

² Intrusion Detection Systems

۴) توابع جداساز غیرخطی را به راحتی و با محاسبات کم، با استفاده از توابع کرنل^۳ تشکیل می‌دهند.

۵) طبقه‌بندهای ماشین بردار پشتیبان در انواع دسته‌بندی‌هایی همچون تشخیص نفوذ به شبکه‌های کامپیوتری، تشخیص ارقام دست‌نویس، تشخیص شیء، شناسایی صورت، دسته‌بندی انواع صداها و مانند آن مورد استفاده قرار گرفته است که در مقایسه با روش‌های دیگر از کارایی قابل ملاحظه‌ای برخوردار است.

از جمله معایب این طبقه‌بند عبارتند از:

- ۱) با افزایش تعداد نمونه‌ها، کارایی این طبقه‌بندها کاهش می‌یابد.
- ۲) با افزایش تعداد بردارهای پشتیبان، نسبت به سایر طبقه‌بندها، در فاز آزمون^۴، کندتر عمل می‌کنند؛ چرا که مرحله آموزش این طبقه‌بندها، معادل با حل یک معادله درجه دوم محدب است که این معادله شامل یک ماتریس نیمه‌معین مثبت با تعداد سطرهایی برابر با تعداد نمونه‌های آموزشی می‌باشد.
- ۳) حساسیت بالا به پارامترها.
- ۴) در صورتی که اندازه مجموعه داده n و تعداد ویژگی‌های آن d باشد، در مرحله آموزش، حداقل پیچیدگی زمانی این گونه طبقه‌بندها $O(n^3 d^3)$ و پیچیدگی حافظه‌ای نیز $O(n^2 d^2)$ خواهد بود. در نتیجه با افزایش نمونه‌ها ماتریس نیز بزرگ‌تر می‌گردد و مسئله برای مجموعه داده‌های حجیم کارایی خود را از دست می‌دهد [Tax '99]. برای رفع این مشکل و تسریع این گونه طبقه‌بندها راهکارهای گوناگونی ارائه شده است.

³ kernel

⁴ Test phase