



دانشکده مهندسی

پایان نامه کارشناسی ارشد در رشته مهندسی کامپیوتر (هوش مصنوعی)

بهبود دقت طبقه بندی با استفاده از روش های انتخاب خصیصه ی ترکیبی

توسط:

فاطمه علیمردانی

استاد راهنما:

دکتر رضا بوستانی

شهریور ۱۳۸۸

به نام خدا

اظہار نامہ

اینجانب ناظم علمیران (۸۵۰۹۸۷) دانشجوی رشته‌ی
مهندسی کامپیوتر گرایش هوش مصنوعی دانشکده‌ی برق و کامپیوتر
اظہار می‌کنم که این پایان نامه حاصل پژوهش خودم بوده و در جاهایی که
از منابع دیگران استفاده کرده‌ام، نشانی دقیق و مشخصات کامل آن را
نوشته‌ام. همچنین اظہار می‌کنم که تحقیق و موضوع پایان نامه‌ام تکراری
نیست و تعهد می‌نمایم که بدون مجوز دانشگاه دستاوردهای آن را منتشر
ننموده و یا در اختیار غیر قرار ندهم. کلیه حقوق این اثر مطابق با آیین‌نامه
مالکیت فکری و معنوی متعلق به دانشگاه شیراز است.

نام و نام خانوادگی: ناظم علمیران
تاریخ و امضا: ۱۸/۶/۱۸

به نام خداوند جان و خرد

به نام خدا

بهبود دقت طبقه بندی با استفاده از روش های انتخاب خصیصه ی
ترکیبی

به وسیله ی :

فاطمه علیمردانی

پایان نامه

ارائه شده به معاونت تحصیلات تکمیلی دانشگاه به عنوان بخشی از
فعالیت های تحصیلی لازم برای اخذ درجه کارشناسی ارشد

در رشته ی:

مهندسی کامپیوتر - گرایش هوش مصنوعی

از دانشگاه شیراز

شیراز

جمهوری اسلامی ایران

ارزیابی شده توسط کمیته پایان نامه با درجه: عالی

دکتر رضا بوستانی، استادیار بخش مهندسی کامپیوتر (رئیس کمیته)

دکتر منصور ذوالقدری جهرمی دانشیار بخش مهندسی کامپیوتر

دکتر کریم رستگار استادیار دانشکده پزشکی دانشگاه علوم پزشکی شیراز

دکتر احمد غنی زاده، دانشیار دانشکده پزشکی دانشگاه علوم پزشکی شیراز

دکتر احمد غنی زاده

کمیته ی برابری جنسیتی - برتر گسالت
فوق تخصص اعصاب و روان کودکان و نوجوانان
دانشگاه شیراز و دانشگاه علوم پزشکی شیراز ن ۳۸۳۳۲۰۰

شهریور ۱۳۸۸

به:

آنهایی که به من آموختند حتی یک کلمه.

سپاسگزاری

در پایان این مسیر پرفراز و نشیب، که کوله‌باری از تجربه‌های گران‌بها را برایم به ارمغان آورده است، به طور خاص از استاد راهنمای محترم آقای دکتر رضا بوستانی، به خاطر صبوری‌ها و راهنمایی‌های ارزنده‌شان سپاسگزارم. همچنین از اساتید بزرگوام جناب آقای دکتر کریم رستگار، جناب آقای دکتر احمد غنی زاده و آقای دکتر منصور ذوالقدر جهرمی که در طول این مدت از حمایت‌های ایشان برخوردار بوده‌ام کمال تشکر را دارم.

از تمامی دوستان به خصوص خانم خدیجه سادات نژاد و کلیه کادر بخش مهندسی و علوم کامپیوتر دانشگاه شیراز که با اینجانب همکاری صادقانه داشتند، تشکر می‌کنم و یک سپاس ویژه از بیماران محترمی که در انجام آزمایشات این تحقیق همکاری بی‌دریغی داشته‌اند دارم و برای تک تک آنها سلامت و شادی روز افزون آرزومندم.

چکیده

بهبود دقت طبقه بندی با استفاده از روش‌های نوین ترکیبی انتخاب خصیصه

به وسیله‌ی:

فاطمه علیمردانی

در سال‌های اخیر، کاراتر سازی روش‌های انتخاب خصیصه به صورت رو به رشد مورد توجه بوده‌اند. تحقیقات انجام شده به منظور کاهش هزینه‌ی محاسبات و همچنین کاهش ریسک تطابق بیش از حد مطالعه شده‌اند. در راستای کاهش بعد خصیصه‌ها از بین همه‌ی الگوریتم‌ها روش‌های پیشرونده به دلیل هزینه‌ی محاسباتی کم بسیار کارا هستند. در این رساله دو الگوریتم پیشرونده‌ی جدید بر اساس مقدار افزایش اطلاعاتی که با حضور هر خصیصه خواهیم داشت ارائه شده است. این الگوریتم‌ها سعی دارند خصیصه‌های حاوی بیشترین اطلاعات را یافته و انتخاب کنند. از طرف دیگر با غلبه بر محدودیت‌های الگوریتم‌های ترتیبی پیشرونده و پسرونده سعی می‌کنند زیرمجموعه‌ای را انتخاب کنند که خصیصه‌های آن بیشترین استقلال را داشته باشند. برای بررسی میزان خوبی خصیصه‌هایی که انتخاب می‌شوند از طبقه بندی کننده‌ی نزدیکترین همسایه استفاده شده است. این طبقه بند وابسته به معیار فاصله‌ای که برای یافتن نزدیکترین همسایه به کار می‌برد می‌باشد. در همین راستا یک نسخه‌ی تصحیحی برای حذف این وابستگی پیشنهاد شده است که از رای گیری بین سه فاصله نتیجه‌ی خود را برای نمونه‌ی آزمایشی اعلام می‌کند. نتایج، بهبود قابل توجه صحت طبقه بندی نسبت به روش‌های شناخته شده را روی داده‌های UCI نشان می‌دهند. دلیل این بهبود، کارایی الگوریتم‌های جدید معرفی شده است و از سویی دیگر طبقه بند پیشنهادی تصمیم پایدارتر و ثابت تری در انتخاب کلاس نمونه‌ی آموزشی نسبت به نسخه‌ی استاندارد آن می‌گیرد.

در یک آزمایش دیگر روش‌های پیشنهادی در تشخیص دو دسته بیماری روانی (دو قطبی و اسکیزوفرنی) که در پزشکی تشخیص کمی روی آن‌ها وجود ندارد به کار برده شده است. این مطالعه آماری روی این دو بیماری برای اولین بار ارائه شده است.

فهرست مطالب

صفحه	عنوان
۱.....	۱- مقدمه
۱	۱-۱- انتخاب خصیصه
۳	۲-۱- ارزیابی یک زیر مجموعه از خصیصه‌ها
۳	۱-۲-۱- معیارهای مستقل
۴	۲-۲-۱- معیارهای وابستگی
۹.....	۲- مروری بر روشهای استخراج خصیصه از سیگنال
۱۱	۱-۲- سیگنالهای EEG
۱۲	۲-۲- تست پاسخ محرک بینایی
۱۲	۳-۲- ثبت سیگنال
۱۸	۴-۲- روشهای استخراج خصیصه
۱۹	۱-۴-۲- مدل خودبازگشتی
۲۱	۲-۴-۲- کورتوسیس
۲۲	۳-۴-۲- بعد فراکتالی و هندسه فراکتالی
۲۳	۱-۳-۴-۲- مولفههای لیپانوف
۲۷	۵-۲- فیلترینگ
۲۸.....	۳- مروری بر روشهای گسسته‌سازی
۲۸	۱-۳- مقدمه
۲۸	۲-۳- معرفی گسسته‌سازی
۲۹	۳-۳- روشهای موجود
۳۱	۴-۳- ارزیابی نتایج گسسته‌سازی
۳۲	۵-۳- چهارچوب گسسته‌سازی
۳۳	۶-۳- روشهای تقسیم کردن
۳۴	۱-۶-۳- بخش‌بندی
۳۷	۲-۶-۳- معیار انتروپی

۴۰	۳-۶-۳- معیار وابستگی
۴۱	۳-۶-۴- معیار دقت
۴۲	۳-۷-۷- روشهای ادغام کردن
۴۲	۳-۷-۱- معیار مربع کای
۴۳	۳-۷-۲- روش Chi Merge
۴۳	۳-۷-۳- روش Chi^2
۴۵	۳-۸-۸- تکنیک گسسته‌سازی بر اساس معیار فیشر
۴۵	۳-۸-۱- معرفی معیار فیشر
۴۶	۳-۸-۲- الگوریتم گسسته‌سازی بر اساس فیشر
۴۸	۴- مروری بر روشهای انتخاب خصیصه
۴۸	۴-۱- مقدمه
۴۹	۴-۲- الگوریتم جستجوی تبو
۵۳	۴-۳- پلاس. ال. ماینس. آر
۵۳	۴-۴- اف. سی. بی. اف
۵۴	۴-۴-۱- ارتباط خطی
۵۵	۴-۴-۲- اطلاعات متقابل
۵۹	۴-۵- شاخص دیویس بولدین
۶۳	۴-۶- رتبه بندی متغیرها
۶۳	۴-۶-۱- اصول متد و علامت گذاری ها
۶۳	۴-۶-۲- شرایط همبستگی
۶۴	۴-۷- مثالهای کوچک اما راه گشا
۶۵	۴-۷-۱- آیا متغیرهای اضافی فرضی به همدیگر کمک میکنند؟
۶۶	۴-۷-۲- چگونه همبستگی بر افزونگی متغیرها تاثیر می گذارد؟
۶۷	۴-۷-۳- آیا متغیری که به تنهایی مفید نیست در کنار بقیه متغیرها سودمند است؟
۶۹	۴-۸-۸- روشهای انتخاب خصیصه‌ی توکار(تو در تو)
۶۹	۴-۸-۱- متدهای زیرمجموعه‌های تودرتو
۷۰	۴-۸-۲- بهینه‌سازی اهداف مستقیم
۷۱	۴-۹-۹- مقولات پیشرفته و مسایل باز
۷۱	۴-۹-۱- واریانس زیرمجموعه‌های انتخاب شده
۷۱	۴-۹-۲- معیارانتخاب خصیصه در مسائل بدون ناظر
۷۲	۴-۹-۳- انتخاب پیشرو در مقابل انتخاب عقبگرد
۷۳	۴-۱۰- مسایل چندکلاسی

۷۴ طبقه بندی کننده	۵-۵
۷۴		۱-۵-۱- مقدمه
۷۵		۲-۵-۲- طبقه بندی بیز ساده
۷۶		۳-۵-۳- طبقه بندی نزدیکترین همسایه
۸۱ روشهای پیشنهادی و نتایج تجربی	۶-۶
۸۱		۱-۶-۱- مقدمه
۸۱		۲-۶-۲- روشهای پیشنهادی
۸۱		۳-۶-۱-۲- الگوریتم حریمانه بر اساس ارتباط کلی خصیصهها
۸۶		۳-۶-۲-۲- الگوریتم سریع کاهش ویژگی با استفاده از شاخص دیویس بولدین
۸۸		۳-۶-۳- مجموعه داده‌های سیگنال
۹۰		۳-۶-۱-۳- روش شکاف آماری (گپ استاتیک)
۹۱		۳-۶-۲-۳- روش ارزیابی Silhouette
۹۴		۳-۶-۳-۳- خوشه بندی تفاضلی
۹۹		۴-۶-۴- داده‌های استاندارد
۱۰۴ نتیجه‌گیری و کارهای آینده	۷-۷
۱۰۴		۱-۷-۱- نتیجه‌گیری
۱۰۴		۲-۷-۲- کارهای آینده
۱۰۶ مراجع	۸-۸

فهرست جداول

صفحه	عنوان
۳۵	جدول 1- نقاط گسستگی داده‌ی با روش پهنای مساوی
۳۵	جدول 2- نقاط گسستگی داده‌ی Iris با روش فراوانی مساوی
۳۶	جدول ۳- نقاط گسستگی در روش 1R
۳۷	جدول 4- مقایسه سه روش بخش بندی
۴۴	جدول 5- مقایسه نتایج سه روش با آزمون مربع کای
۹۰	جدول 6 - تعداد خصیصه‌های انتخاب شده توسط روشهای مختلف.
۹۵	جدول ۷- دقت نتیجه شده از NN استاندارد در مقایسه با NN پیشنهادی
۹۷	جدول ۸- دقت طبقه بندی پیشنهادی بر روی هر کدام از خصیصه‌ها به تنهایی
۹۷	جدول ۹- نتایج انتخاب خصیصه NN+ استاندارد و NN معرفی شده روی داده‌های نویزی
۱۰۰	جدول ۱۰- - مشخصات مجموعه داده‌های UCI به کار رفته در این رساله
۱۰۱	جدول ۱۱- - نتایج طبقه بندی کننده‌ی پیشنهادی روی خصیصه‌های انتخاب شده مختلف
۱۰۳	جدول ۱۲- - نتایج طبقه بندی کننده‌های مختلف روی خصیصه‌های انتخابی GOR

فهرست شکل‌ها

صفحه	عنوان
۵	شکل ۱- دسته بندی روشهای انتخاب خصیصه [1]
۱۳	شکل ۲- مکان الکترودها در سیستم ۲۰-۱۰
۱۵	شکل ۳- الف) نمونه‌های از سیگنالهای ضبط شده الف) دوقطبی، ب) اسکیزوفرنی
۱۷	شکل ۴- طیف توان نمونه‌های از هر دو نوع بیماری به کمک Color Map، الف) دو قطبی، ب) اسکیزوفرنی
۲۰	شکل ۵- مدلسازی AR [77]
۲۰	شکل ۶- یک نمونه سیستم خطی [76]
۲۲	شکل ۷- الف) یک عملگر ساده فرکتال و اثر متوالی آن بر یک پاره خط. ب) یک شکل- فرکتال مصنوعی [80].
۲۶	شکل ۸- منحنی جمع همبستگی برحسب شعاع کره برای یک سیگنال غیر نویزی.
۳۳	شکل ۹- چهارچوب سلسله مراتبی گسسته سازی [25]
۳۶	شکل ۱۰- مثال عملی مراحل انجام روش IR
۴۷	شکل ۱۱- فلوچارت تکنیک گسسته سازی با معیار فیشر [74]
۴۸	شکل ۱۲- چهار گام کلیدی در فرایند انتخاب خصیصه [21]
۵۱	شکل ۱۳- شمای کلی الگوریتم تبو [15]
۵۲	شکل ۱۴- شبه کد الگوریتم تبو [20]
۵۴	شکل ۱۵- شبه کد الگوریتم پلاس. ال. ماینس. آر [83]
۵۹	شکل 16- شبه کد الگوریتم اف. سی. بی. اف [10]
۶۱	شکل ۱۷- کارایی روشهای مختلف انتخاب خصیصه تحت شرایط خاص [5]
۶۵	شکل ۱۸- اطلاعات به دست آمده از متغیرهای اضافی فرضی
۶۶	شکل ۱۹- تاثیر همبستگی بر افزونگی داده‌ها
۶۹	شکل ۲۰- بررسی مفید بودن متغیرها در همکاری با هم.
۷۳	شکل ۲۱- بررسی انتخاب عقبگرد و پیشرو
۷۷	شکل ۲۲- مثالی از عملکرد NN با $k=5$
۷۹	شکل ۲۳- کانتورهایی از فاصله‌های مساوی [28]
۸۶	شکل ۲۴- الگوریتم حریرصانه بر اساس ارتباط کلی خصیصه‌ها.

- شکل ۲۵- شبه کد الگوریتم DB-FFR ۸۸
- شکل ۲۶- مقدار Silhouett برای هر دو دسته بیماری ۹۳
- شکل ۲۷- مقایسه‌ی نرخ تشخیص بین دو بیماری با الگوریتم های مختلف انتخاب ۹۶
- شکل ۲۸- دقت روشهای مختلف انتخاب + NN معرفی شده روی داده‌های نویزی ۹۸
- شکل ۲۹- مقایسه روشهای مختلف انتخاب خصیصه با طبقه بند پیشنهادی ۱۰۲
- شکل ۳۰- مقایسه نتیجه‌ی GOR با استفاده از طبقه بند های مختلف ۱۰۳

فهرست اختصارات

ABBREVIATION	MEANING
AR	Auto Regressive
BMD	Bipolar Mode Disorder
DB	Davies Boildin
DB-FFR	Davies Bouldin Fast Feature Reduction
GOR	Greedy Overall Relevancy
MNN	Modified Nearest Neighbor
ONR	Overall Normalize Relevancy
R-N-R	Relevant-Non- Redundant

فصل اول

مقدمه

۱- مقدمه

۱-۱- انتخاب خصیصه

جمع آوری داده با انجام آزمایش کار بسیار ساده و راحتی است. داده ها با سرعت بی سابقه ای گرد آوری و ذخیره می شوند که بسیار بیشتر از سرعت پردازش انسان است بنابراین پیش پردازش این داده های تل انبار شده به منظور استخراج یک مجموعه داده کارا تر امر بسیار مهمی در مسائل یادگیری ماشین و داده کاوی به شمار می رود. انتخاب خصیصه از رایج ترین تکنیکهاست که کارایی آن در برخورد با مجموعه داده های بزرگ در بسیاری از مسائل یادگیری ماشین از قبیل: طبقه بندی، خوشه سازی و داده کاوی از سال ۱۹۷۰ ثابت شده است [1,2,3,4]. رشهای انتخاب خصیصه در دو دسته کلی یادگیری ماشین یعنی روشهای همره با ناظر و نیز روشهای بدون ناظر به کار برده شده اند اما در این رساله ما توجه خود را به مسائل همراه با ناظر (طبقه بندی) متمرکز کرده ایم که در آنها برچسب کلاس هر کدام از داده ها از ابتدا در مسئله مشخص است. پس یک مرور کلی بر روشهای انتخاب خصیصه در مسائل طبقه بندی کننده در این رساله انجام می شود. برای هر کدام از روشها مهمترین زمینه های کاربردی بازنگری خواهد شد و در پایان یک روش جدید ارائه خواهد شد تا در حالیکه سعی دارد فواید روشهای قبلی را به ارث ببرد با دقت بهتری خصیصه های بهینه را گزینش کند.

انتخاب یک زیر مجموعه از خصیصه ها به دو موضوع برای بررسی نیازمند است: ۱- یک الگوریتم و استراتژی جستجو کننده برای انتخاب خصیصه های کاندید و ۲- یک تابع هدف برای ارزیابی خصیصه های کاندید تا میزان خوبی این کاندیدها مشخص شود. از طرفی باید نشانه ای از مقدار ارزیابی شده توسط تابع هدف، به عنوان بازخورد به الگوریتم جستجو کننده باز گردانده شود تا در انتخاب کاندیدهای بعدی در نظر گرفته شود. الگوریتم های جستجوگر به طور کلی به سه دسته تقسیم می شوند: الف) الگوریتم های نمایی^۱، ب) الگوریتم های ترتیبی^۲ و الگوریتم های تصادفی^۳. و از طرفی الگوریتم های ترتیبی خود به دو گروه روش های پیشرونده^۴

¹ exponential Algorithms

² Sequential Algorithms

³ Randomized Algorithms

⁴ Forward

و پسرونده⁵ دسته بندی می‌شوند. روش‌های پسرونده کار خود را با مجموعه‌ای با حضور تمام خصیصه‌ها شروع می‌کند و سپس به طور ترتیبی خصیصه‌هایی که بیشترین کاهش در تابع هدف را موجب می‌شوند از این مجموعه حذف می‌کند. و روش‌های پیشرونده ابتدا مجموعه جواب را خالی فرض می‌کند و گام به گام در هر مرحله خصیصه‌ای را به جواب اضافه می‌کند که بیشترین تغییر را در مقدار تابع هدف موجب می‌شود.

بطور کلی توابع هدف مورد استفاده در روش‌های انتخاب خصیصه در دو دسته قرار می‌گیرند [32]: فیلترها⁶ و رپرها⁷. در رویکرد فیلتر، انتخاب بهترین زیرمجموعه از خصیصه‌ها در واقع یک مرحله پیش‌پردازش است و بعداً خصیصه‌ها برای طبقه بندی به الگوریتم‌های یادگیری ماشین داده می‌شوند. این روش بر اساس یک معیار از پیش تعیین شده می‌باشد که مستقل از کارایی خصیصه‌های انتخاب شده در افزایش دقت طبقه بند استفاده شده است. هر چند رویکرد فیلتر زمان اجرای کمتری لازم دارد اما این عیب را دارد که به طبقه بندی کننده مورد استفاده در فاز بعدی مسئله توجه نمی‌کند. در حالیکه در رویکرد رپر، همکنش بین مجموعه خصیصه انتخاب شده و طبقه بندی کننده به طور خاص مورد توجه روش جستجو در خصیصه‌ها قرار می‌گیرد. به عبارتی رپرها یک ماشین یادگیری خاص برای جستجو زیرمجموعه خصیصه مناسب بکار می‌برند. ماشین یادگیری روی یکی از زیرمجموعه خصیصه‌های مورد نظر آموزش داده می‌شود و درجه دقت یادگیرنده در طبقه بندی داده‌ها به عنوان معیار ارزیابی زیرمجموعه ای از خصیصه‌های کاندید بازگردانده می‌شود. این مرحله برای همه‌ی زیرمجموعه‌های کاندید از خصیصه‌ها تکرار می‌شود تا جایی که بهترین زیرمجموعه خصیصه یافت شود. رویکرد رپر توسط John و Kohavi معرفی شد، و بطور کلی می‌تواند طبقه بندی کننده‌ای نهایی با دقت پیشگویی خوب تولید کند [32]؛ اما رویکرد رپر در مقایسه با روش فیلتر هزینه‌ی محاسباتی بیشتری دارد و احتمال داشتن مشکل انطباق بیش از حد⁸ را افزایش می‌دهد [33,34]. بنابراین در بسیاری کاربردها، روش‌های فیلتر نسبت به رپر ترجیح داده می‌شوند.

معیار انتخاب خصیصه نقش حیاتی در روش‌های فیلتر بازی می‌کند. در مبحث طبقه بندی، این معیار سعی در اندازه گیری توانایی یک خصیصه یا یک مجموعه خصیصه برای جدا کردن کلاس‌های مختلف دارد. تا به حال معیارهای مختلفی مانند معیارهای فاصله، وابستگی، و سازگاری برای انتخاب خصیصه بکار رفته است [35,36-37]. ولی این معیارها ممکن است به مقادیر قطعی داده‌های آموزشی بسیار حساس باشند؛ در نتیجه آنها به راحتی تحت تاثیر نویز⁹ یا داده‌های بیرونی¹⁰ قرار می‌گیرند. در حالیکه معیارهای اطلاعات، از قبیل انترپی¹¹ یا

⁵ Backward

⁶ Filter approach

⁷ wrapper

⁸ Over-fitting

⁹ Noise

¹⁰ Outliers

اطلاعات متقابل، میزان اطلاعات یا عدم قطعیت یک خصیصه را برای طبقه‌بندی بررسی می‌کنند. بر تئوری اطلاعات [38, 39] Shannon، مدلی که مقدار اطلاعات را افزایش ندهد بی‌فایده تلقی می‌شود و انتظار نمی‌رود دقت پیشگویی آن بهتر از تنها یک حدس تصادفی باشد [40]. یکی از محاسن معیار میزان اطلاعات یک خصیصه این است که این معیار تنها به توزیع احتمال یک متغیر تصادفی وابسته است، نه به مقادیر قطعی آن. معیارهای اطلاعات بطور وسیع برای انتخاب خصیصه مورد استفاده قرار گرفته‌اند [41, 42]، از قبیل بسیاری الگوریتم‌های یادگیری مانند [43] C4.5. یکی از این معیارها اطلاعات متقابل^{۱۲} می‌باشد.

۱-۲- ارزیابی یک زیر مجموعه از خصیصه‌ها

در همه کاربردها، به منظور ارزیابی یک زیر مجموعه از خصیصه‌ها میزان خوبی این زیر مجموعه توسط یک معیار معین به دست می‌آید. معیارهای ارزیابی به طور وسیع با توجه به وابستگی آنها به الگوریتم کاوشی که بعداً روی زیرمجموعه انتخاب شده عمل خواهد کرد در دو گروه دسته بندی می‌شوند: معیارهای مستقل^{۱۳} و معیارهای وابسته^{۱۴}.

۱-۲-۱- معیارهای مستقل

معیارهای مستقل عموماً در الگوریتم‌های فیلتر استفاده می‌شوند. این معیارها سعی دارند میزان خوبی یک خصیصه را با بهره‌گیری از ویژگیهای ذاتی داده‌های آموزشی و بدون در نظر گرفتن الگوریتم‌های کاوش اندازه‌گیری کنند.

برخی معیارهای مستقل عبارتند از مقیاس فاصله، مقدار اطلاعات، مقدار همبستگی و میزان سازگاری یا انطباق [5,6]. مقیاس فاصله با نامهای دیگری مثل مقدار جدایی، انحراف و کجی یا مقدار تشخیص پذیری نیز به کار برده شده است. از آنجا که در یک مسئله دو کلاسه هدف یافتن خصیصه‌ای است که کلاس‌ها را تا حد امکان از هم جدا کند، خصیصه X نسبت به خصیصه Y ارجح است اگر با در نظر گرفتن خصیصه X فاصله بیشتری بین احتمال شرطی کلاس‌ها ایجاد شود. یک خصیصه را غیر قابل تمیز دادن می‌گوییم اگر فاصله بین کلاس‌ها در حضور آن صفر شود. معیار مقدار اطلاعات معمولاً میزان افزایش اطلاعات را معین می‌کنند. میزان افزایش تفاوت بین عدم قطعیت از پیش مشخص شده در داده‌ها وقتی خصیصه X را

¹¹ Entropy

¹² Mutual information

¹³ Independent Criteria

¹⁴ Dependency Criteria

استفاده می‌کنیم و عدم قطعیت پسین با حضور خصیصه X است. هرچه مقدار اطلاعات یک خصیصه بیشتر باشد مطلوب تر است پس در هر مرحله خصیصه با مقدار اطلاعات بیشتر انتخاب می‌شود.

معیارهایی نیز به عنوان اندازه‌ای از همبستگی و شباهت شناخته شده‌اند و در واقع توانایی پیش‌بینی کردن یک مقدار از مقادیر دیگر را اندازه می‌گیرند. در مسئله انتخاب خصیصه برای طبقه بندی ما به دنبال خصیصه‌هایی هستیم که بیشترین همبستگی را با برچسب کلاس هر داده دارد. و در مسئله انتخاب خصیصه برای خوشه‌سازی داده‌ها، میزان شباهت بین دو خصیصه مد نظر قرار می‌گیرد.

معیارهای سازگاری مشخصه‌های متفاوتی نسبت به معیارهای ذکر شده دارند. این معیارها به داده‌های کلاس متکی هستند [6] و سعی می‌کنند مینیمم تعداد خصیصه‌ای را پیدا کنند که مطابق با وقتی که همه خصیصه‌ها را در نظر می‌گیریم بتواند کلاس‌ها را از هم جدا کند. عدم انطباق دو خصیصه به معنی این است که علیرغم داشتن مقدار یکسان در این خصیصه‌ها دو برچسب متفاوت بر اساس آنها انتخاب شود.

۱-۲-۲- معیارهای وابستگی

در رپرها از معیارهای همبستگی استفاده می‌شود و لازم است یک الگوریتم کاوش از قبل تعریف شده داشته باشند. رپرها اغلب کارایی بیشتری دارند چون خصیصه‌هایی که برای این الگوریتم کاوش مناسب ترند را پیدا می‌کنند ولی در عوض هزینه محاسباتی زیادی دارند و در ضمن ممکن است برای الگوریتم‌های کاوش دیگر، کارایی آنها کاهش پیدا کند.

همه‌ی روشهای انتخاب خصیصه گفته شده به یک معیار یا شرط برای توقف کار نیاز دارند. معیار توقف معمولاً یکی از موارد زیر است:

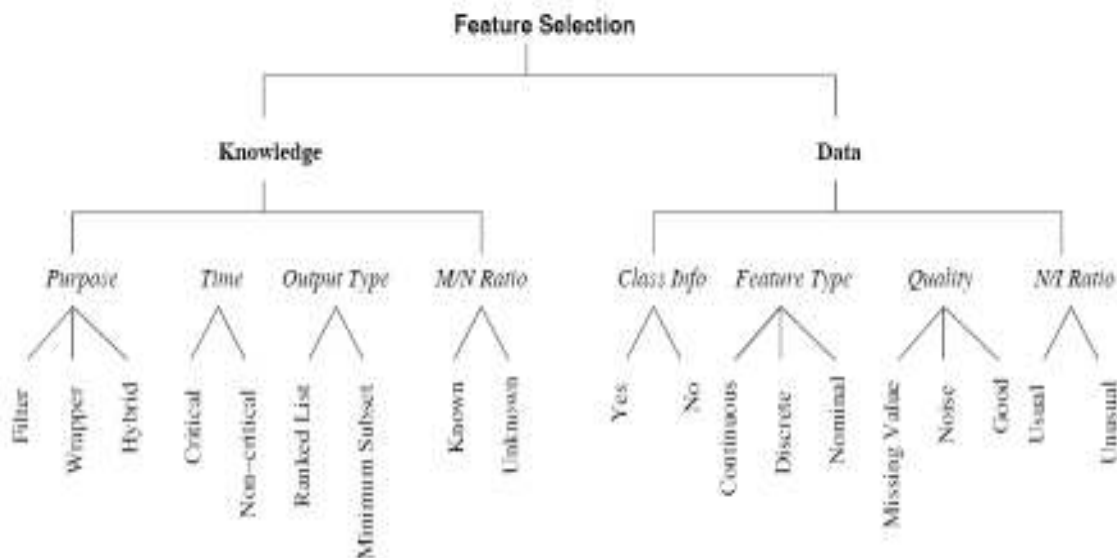
۱- جستجو کامل شده باشد،

۲- یک حد بالایی مشخصی مرتفع شده است (به عنوان مثال به مینیمم تعداد خصیصه‌ها رسیده‌ایم و یا به ماکزیمم تعداد تکرار) ،

۳- دیگر اضافه یا حذف کردن خصیصه‌ی بعدی تاثیری بر کارایی مجموعه خصیصه انتخاب شده تاکنون ندارد و

۴- وقتی که یک زیر مجموعه‌ی به اندازه کافی خوب از خصیصه‌ها انتخاب شده باشد. در نهایت پس از انتخاب زیر مجموعه‌ی از خصیصه‌ها با یکی از روشهای مورد نظر وقت آن است که این مجموعه را ارزیابی کنیم. یک راه ساده برای ارزیابی نتایج آن است که مستقیماً

با استفاده از اطلاعاتی که از قبل راجع به داده‌ها داریم نتیجه را مقدار سنجی کنیم. اگر از قبل خصیصه‌های مرتبط و مطرح در جواب بهینه را بدانیم می‌توانیم این مجموعه معین را با مجموعه انتخاب شده مقایسه کنیم. اینکه اطلاعاتی درباره خصیصه‌های غیر مرتبط و افزونه داشته باشیم هم می‌تواند مفید باشد. زیرا انتظار داریم این خصیصه‌ها انتخاب نشوند. البته در کاربردهای دنیای واقعی معمولاً این دانش را از قبل راجع به خصیصه‌ها نداریم بنابراین باید کار را به طور غیر مستقیم دنبال کنیم، به این ترتیب که بر میزان تغییر در کارایی الگوریتم کاوش (طبقه بند یا خوشه ساز) وقتی خصیصه‌های شرکت کننده را تغییر می‌دهیم نظارت کنیم. مثلاً اگر نرخ خطای طبقه بندی را به عنوان مشخص کننده کارایی عمل کاوش انتخاب کنیم می‌توانیم به سادگی این نرخ را قبل و بعد از اضافه یا حذف یک خصیصه حساب کنیم. [5]. در نمودار ۱ یک قالب دسته بندی از روشهای انتخاب خصیصه آمده است.



شکل ۱ - دسته بندی روش‌های انتخاب خصیصه [1]

همان‌طور که می‌بینید دو فاکتور اصلی دانش از قبل مشخص شده و داده‌ها هستند. دانش خود شامل: هدف، زمان مورد نظر، نوع خروجی و نرخ $\frac{M}{N}$ است که در آن M تعداد مورد انتظار برای زیر مجموعه ای که می‌خواهیم انتخاب کنیم و N تعداد کل مجموعه اصلی خصیصه-