

به نام یزدان پاک



دانشگاه کردستان
دانشکده فنی و مهندسی
گروه مهندسی کامپیوتر و فناوری اطلاعات

عنوان:

راهکار ترکیبی برای انتخاب ویژگی در داده‌های ابعاد بالا

پژوهشگر:

محمدحسین دشتبان

استاد راهنما:

دکتر پرهام مرادی

استاد مشاور:

دکتر هادی زارع

پایان نامه کارشناسی ارشد رشته مهندسی کامپیوتر گرایش هوش مصنوعی

اسفند ماه ۱۳۹۱

کلیه حقوق مادی و معنوی مترتب بر نتایج مطالعات،

ابتکارات و نوآوری های ناشی از تحقیق موضوع

این پایان نامه (رساله) متعلق به دانشگاه کردستان است.

*** تعهد نامه ***

اینجانب محمدحسین دشتبان دانشجوی کارشناسی ارشد رشته مهندسی کامپیوتر گرایش هوش مصنوعی دانشگاه کردستان، دانشکده فنی و مهندسی گروه مهندسی کامپیوتر و فناوری اطلاعات تعهد می نمایم که محتوای این پایان نامه نتیجه تلاش و تحقیقات خود بوده و از جایی کپی برداری نشده و به پایان رسانیدن آن نتیجه تلاش و مطالعات مستمر اینجانب و راهنمایی استاد بزرگوارم بوده است.

با تقدیم احترام

محمدحسین دشتبان

۱۳۹۱/۱۲/۰۱

تقدیم با احترام

به قهرمان جاویدان ایران

کورش کبیر



دانشگاه کردستان

دانشکده فنی و مهندسی

گروه مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی ارشد رشته مهندسی کامپیوتر گرایش هوش مصنوعی

عنوان:

راهکار ترکیبی برای انتخاب ویژگی در داده‌های ابعاد بالا

پژوهشگر:

محمدحسین دشتیان

در تاریخ..... توسط کمیته تخصصی و هیات داوران زیر مورد بررسی قرار گرفت و با نمره..... و درجه..... به تصویب رسید.

<u>امضاء</u>	<u>مرتبۀ علمی</u>	<u>نام و نام خانوادگی</u>	<u>هیات داوران</u>
	استادیار	دکتر پرهام مرادی	۱- استاد راهنما
	استادیار	دکتر هادی زارع	۲- استاد مشاور
	استادیار	دکتر حمیدفرورش	۳- استاد داور خارجی
	استادیار	دکتر فردین اخلاقیان	۴- استاد داور داخلی

مهر و امضاء معاون آموزشی و تحصیلات

مهر و امضاء گروه

تکمیلی دانشکده

چکیده

با پیشرفت روزافزون تکنولوژی در زمینه داده‌کاوی در حوزه‌های علمی مختلف، مجموعه‌داده‌های با ابعاد بسیار بالا در حال افزایش است که منجر به کاهش کارایی الگوریتم‌های دسته‌بندی می‌شود. لذا نیاز به کاهش حجم این مجموعه داده‌ها امری ضروری است. در مجموعه‌داده‌ها با ابعاد بالا، تعداد زیادی ویژگی برای هر نمونه وجود دارد که بسیاری از آنها نامرتب و زاید می‌باشند. در این پایان‌نامه بر روی انتخاب ویژگی بر روی مجموعه‌داده‌های ابعاد بالای دو حوزه مختلف علم، بیوانفورماتیک و متن، کار شده است. برای هر یک از این حوزه‌ها راهکارهای انتخاب ویژگی متفاوتی توسط محققان ارائه شده است که این راهکارها وابسته به ماهیت ویژگی‌های حوزه مورد نظر می‌باشد. مثلاً ویژگی‌های داده‌های میکروآرایه مقدار "بیان ژن‌ها" می‌باشند که عددی حقیقی می‌باشد در حالی که در متن، ویژگی‌ها واژه‌ها بوده که الگوریتم‌های ارائه شده در این حوزه بر روی خصوصیت آماری آنها که ماهیتی گسسته دارد تمرکز دارد.

راهکارهای ارائه شده برای انتخاب ویژگی به دو دسته کلی باناظر و بی‌ناظر تقسیم بندی می‌شوند. راهکارهای باناظر از برچسب کلاس‌ها در انتخاب ویژگی کمک می‌گیرند، در حالی که در حالت بی‌ناظر تنها از مقادیر ویژگی‌ها استفاده می‌شود. تحلیل واریانس از راهکارهای بی‌ناظر می‌باشد که از دیرباز مورد توجه محققان بوده است. در قسمت اول این پایان‌نامه، روش‌های انتخاب ویژگی بی‌ناظر و باناظر با تکیه بر استخراج ویژگی، تحلیل واریانس و خوشه‌بندی پیشنهاد شده است. روش ارائه شده بر روی شش مجموعه داده بزرگ بیوانفورماتیک که ویژگی‌های آن ژن‌ها می‌باشند، اعمال شده است. آزمایشات و بررسی‌های مختلف انجام گرفته نشان می‌دهند که روش بی‌ناظر و باناظر پیشنهادی در مجموعه داده‌های مختلف کارایی قابل قبولی را کسب نموده است.

در راهکار پیشنهادی دوم پایان‌نامه، روش انتخاب ویژگی مبتنی بر فیلتر با تکیه بر عامل‌های احتمالاتی تاثیرگذار در دسته‌بندی متن که در روش‌های انتخاب ویژگی احتمالاتی پرکاربرد به کار رفته، ارائه می‌شود. روش ارائه شده از جنبه‌های مختلف مورد تحلیل قرار گرفته و کارایی ویژگی‌های انتخابی آن در دسته‌بندی متن با روش‌های دیگر انتخاب ویژگی مبتنی بر فیلتر مقایسه شده است. آزمایشات متعدد، روش‌های فیلتر را از جنبه‌های مختلف همانند: میزان اشتراک ویژگی‌های برتر انتخاب شده، بررسی واریانس ویژگی‌ها، کارایی ویژگی‌های انتخاب شده بر اساس معیارهای مختلف، رفتار کارایی آنها با افزایش تعداد ویژگی‌ها و میزان دقت و بازیابی روش‌ها نسبت به یکدیگر، به طور عملی مورد مطالعه قرار می‌دهند. سه مجموعه داده استاندارد: Reuter-R8، 20Newsgroup و WebKB در این مطالعه استفاده شده است. آزمایشات مختلف نشان دهنده این است که روش پیشنهادی در هر سه مجموعه داده توانایی رقابت با روش‌های موفق انتخاب ویژگی مبتنی بر فیلتر را داراست به طوریکه در برخی موارد اختلاف قابل توجهی را ایجاد نموده است.

کلمات کلیدی: انتخاب ویژگی، استخراج ویژگی، تحلیل واریانس، طبقه بندی متن، روش های مبتنی بر فیلتر، انتخاب ژن، رتبه‌دهی ژن‌ها، تحلیل میکروآرایه.

فهرست مطالب

صفحه

عنوان

۱	فصل ۱ مقدمه.....
۱	۱-۱ سابقه و انگیزه تحقیق
۳	۲-۱ مسئله‌ی تحقیق و اهداف پایان‌نامه
۴	۳-۱ دورنمای پایان‌نامه.....
۶	فصل ۲ پیشینه.....
۶	۱-۲ مقدمه.....
۶	۲-۲ روش‌های انتخاب ویژگی
۷	۳-۲ انتخاب ویژگی در متن
۸	۳-۲-۱ روش روپوش
۸	۳-۲-۲ انتخاب ویژگی مبتنی بر فیلتر
۹	۳-۲-۳ روش‌های ترکیبی
۱۰	۳-۲-۴ استخراج ویژگی
۱۰	۳-۲-۵ روش‌های ادغام شده.....
۱۱	۴-۲ روش‌های استفاده شده در پایان‌نامه.....
۱۱	۴-۲-۱ بهره‌ی اطلاعاتی
۱۲	۴-۲-۲ روش شاخص جینی
۱۲	۴-۲-۳ روش‌هی دو
۱۳	۴-۲-۴ اطلاعات متقابل
۱۵	۴-۲-۵ روش TFIDF
۱۵	۴-۲-۶ روش DFS
۱۷	۴-۲-۷ روش LapAOFS
۲۱	۴-۲-۸ روش رتبه‌دهی Fisher
۲۳	فصل ۳ انتخاب ویژگی استخراجی.....
۲۳	۱-۳ مقدمه.....
۲۴	۲-۳ فرض‌های اولیه.....
۲۴	۳-۳ انتخاب ویژگی استخراجی LapAOFS

۲۵	۳-۴ انتخاب ویژگی استخراجی باناظر.....
۲۷	۳-۵ مجموعه داده.....
۲۸	۳-۶ نتایج عملی.....
۲۹	۳-۶-۱ ارزیابی کارایی.....
۲۹	۳-۶-۲ نتایج روش انتخاب ویژگی استخراجی LapAOFS.....
۳۳	۳-۶-۲ نتایج روش استخراجی باناظر Fisher.....
۳۸	۳-۶-۳ مقایسه ویژگی‌های انتخاب شده روش استخراجی.....
۴۱	۳-۷ نتیجه‌گیری.....

فصل ۴ انتخاب ویژگی مبتنی بر فیلتر پیشنهادی..... ۴۲

۴۲	۴-۱ مقدمه.....
۴۲	۴-۲ روش فیلتر پیشنهادی.....
۴۳	۴-۲-۱ عامل‌های احتمالاتی.....
۴۴	۴-۲-۲ فیلتر CF.....
۴۵	۴-۲-۳ روش پیشنهادی.....
۴۷	۴-۳ مجموعه داده.....
۴۹	۴-۴ معیارهای ارزیابی.....
۵۰	۴-۵ ارزیابی نتایج.....
۵۰	۴-۵-۱ پیش پردازش.....
۵۱	۴-۵-۲ بررسی واریانس.....
۵۴	۴-۵-۳ ویژگی‌های مشترک روش پیشنهادی.....
۵۶	۴-۵-۴ ویژگی مشترک روش‌ها با هم.....
۵۷	۴-۵-۵ ارزیابی بر روی مجموعه داده Newsgroup.....
۶۱	۴-۵-۶ ارزیابی بر روی مجموعه داده Reuter-R8.....
۶۵	۴-۵-۷ ارزیابی بر روی مجموعه داده WebKB.....
۶۹	۴-۵-۸ انتخاب ویژگی تصادفی.....
۷۰	۶-۴ نتیجه‌گیری.....

فصل ۵ جمع‌بندی..... ۷۲

۷۲	۵-۱ کارهای انجام شده.....
۷۳	۵-۲ کارهای آینده.....
۷۵	فهرست منابع.....

فهرست شکل‌ها

- شکل ۳-۱: مراحل روش انتخاب ویژگی استخراجی LAPAOFS ۲۴
- شکل ۳-۲: مراحل روش استخراجی پیشنهادی باناظر ۲۶
- شکل ۳-۳: روش پیشنهادی بی‌ناظر (دایره) در مقابل روش LAPAOFS بر روی مجموعه داده COLON ۳۰
- شکل ۳-۴: روش پیشنهادی بی‌ناظر (دایره) در مقابل روش LAPAOFS بر روی مجموعه داده SRBCT ۳۰
- شکل ۳-۵: روش پیشنهادی بی‌ناظر (دایره) در مقابل روش LAPAOFS بر روی مجموعه داده LEUKMIA ۳۱
- شکل ۳-۶: روش پیشنهادی بی‌ناظر (دایره) در مقابل روش LAPAOFS بر روی مجموعه داده HEDENFALK ۳۱
- شکل ۳-۷: روش پیشنهادی بی‌ناظر (دایره) در مقابل روش LAPAOFS بر روی مجموعه داده BREAST ۳۲
- شکل ۳-۸: روش پیشنهادی بی‌ناظر (دایره) در مقابل روش LAPAOFS بر روی مجموعه داده GCM ۳۲
- شکل ۳-۹: روش پیشنهادی باناظر (دایره) در مقابل روش FISHER بر روی مجموعه داده COLON ۳۳
- شکل ۳-۱۰: روش پیشنهادی باناظر (دایره) در مقابل روش FISHER بر روی مجموعه داده SRBCT ۳۳
- شکل ۳-۱۱: روش پیشنهادی باناظر (دایره) در مقابل روش FISHER بر روی مجموعه داده LEUKMIA ۳۴
- شکل ۳-۱۲: روش پیشنهادی باناظر (دایره) در مقابل روش FISHER بر روی مجموعه داده HEDENFALK ۳۴
- شکل ۳-۱۳: روش پیشنهادی باناظر (دایره) در مقابل روش FISHER بر روی مجموعه داده BREAST ۳۴
- شکل ۳-۱۴: روش پیشنهادی باناظر (دایره) در مقابل روش FISHER بر روی مجموعه داده GCM ۳۵
- شکل ۳-۱۵: روش پیشنهادی باناظر (دایره آبی) در مقابل خوشه‌بندی بر روی مجموعه داده COLON ۳۶
- شکل ۳-۱۶: روش پیشنهادی باناظر (دایره آبی) در مقابل خوشه‌بندی بر روی مجموعه داده SRBCT ۳۶
- شکل ۳-۱۷: روش پیشنهادی باناظر (دایره آبی) در مقابل خوشه‌بندی بر روی مجموعه داده LEUKMIA ۳۶
- شکل ۳-۱۸: روش پیشنهادی باناظر (دایره آبی) در مقابل خوشه‌بندی بر روی مجموعه داده HEDENFALK ۳۷
- شکل ۳-۱۹: روش پیشنهادی باناظر (دایره آبی) در مقابل خوشه‌بندی بر روی مجموعه داده BREAST ۳۷
- شکل ۳-۲۰: روش پیشنهادی باناظر (دایره آبی) در مقابل خوشه‌بندی بر روی مجموعه داده COLON ۳۷
- شکل ۳-۲۱: ویژگی‌های مشترک روش استخراج ویژگی پیشنهادی و FISHER در ۱۰ ویژگی برتر ۳۸
- شکل ۳-۲۱: ویژگی‌های مشترک روش استخراج ویژگی پیشنهادی و FISHER در ۱۰۰ ویژگی برتر ۳۹
- شکل ۳-۲۲: ویژگی‌های مشترک روش بی‌ناظر پیشنهادی و LAPAOFS در ۱۰ ویژگی برتر ۴۰
- شکل ۳-۲۳: ویژگی‌های مشترک روش بی‌ناظر پیشنهادی و LAPAOFS در ۱۰۰ ویژگی برتر ۴۰
- شکل ۴-۱: مراحل روش فیلتر پیشنهادی برای انتخاب ویژگی در متن ۴۶
- شکل ۴-۲: مراحل پیش‌پردازش مجموعه داده‌های متنی ۵۱
- شکل ۴-۳: بررسی واریانس ویژگی‌ها بر روی مجموعه داده REUTER-R8 ۵۲
- شکل ۴-۴: بررسی واریانس ویژگی‌ها بر روی مجموعه داده WEBKB ۵۲
- شکل ۴-۵: بررسی واریانس ویژگی‌ها بر روی مجموعه داده NEWSGROUP ۵۳
- شکل ۴-۶: میزان اشتراک ویژگی روش فیلتر پیشنهادی بر روی مجموعه داده NEWSGROUP ۵۴
- شکل ۴-۷: میزان اشتراک ویژگی روش فیلتر پیشنهادی بر روی مجموعه داده WEBKB ۵۵
- شکل ۴-۸: میزان اشتراک ویژگی روش فیلتر پیشنهادی بر روی مجموعه داده REUTER-R8 ۵۵
- شکل ۴-۹: کارایی روشهای فیلتر با ۵۰ ویژگی برتر بر روی مجموعه داده NEWSGROUP ۵۸
- شکل ۴-۱۰: کارایی روشهای فیلتر با ۱۰۰ ویژگی برتر بر روی مجموعه داده NEWSGROUP ۵۹
- شکل ۴-۱۱: دقت و بازیابی روشهای فیلتر با ۱۰۰ ویژگی برتر بر روی مجموعه داده NEWSGROUP ۶۰
- شکل ۴-۱۲: کارایی میکرو F1 روشهای فیلتر با افزایش ویژگی در مجموعه داده NEWSGROUP ۶۱
- شکل ۴-۱۳: کارایی ماکرو F1 روشهای فیلتر با افزایش ویژگی در مجموعه داده NEWSGROUP ۶۱
- شکل ۴-۱۴: کارایی روشهای فیلتر با ۵۰ ویژگی برتر بر روی مجموعه داده REUTER-R8 ۶۲
- شکل ۴-۱۵: کارایی روشهای فیلتر با ۱۰۰ ویژگی برتر بر روی مجموعه داده REUTER-R8 ۶۳

- شکل ۴-۱۶: کارایی روشهای فیلتر با ۱۵۰ ویژگی برتر بر روی مجموعه داده REUTER-R8..... ۶۳
- شکل ۴-۱۷: دقت و بازیابی روشهای فیلتر با ۱۵۰ ویژگی برتر بر روی مجموعه داده REUTER-R8..... ۶۴
- شکل ۴-۱۸: کارایی میکرو F1 روشهای فیلتر با افزایش ویژگی در مجموعه داده REUTER-R8..... ۶۵
- شکل ۴-۱۹: کارایی ماکرو F1 روشهای فیلتر با افزایش ویژگی در مجموعه داده REUTER-R8..... ۶۵
- شکل ۴-۲۰: کارایی روشهای فیلتر با ۵۰ ویژگی برتر بر روی مجموعه داده WEBKB..... ۶۶
- شکل ۴-۲۱: کارایی روشهای فیلتر با ۱۰۰ ویژگی برتر بر روی مجموعه داده WEBKB..... ۶۶
- شکل ۴-۲۲: کارایی روشهای فیلتر با ۱۵۰ ویژگی برتر بر روی مجموعه داده WEBKB..... ۶۷
- شکل ۴-۲۳: دقت و بازیابی روشهای فیلتر با ۵۰ ویژگی برتر بر روی مجموعه داده WEBKB..... ۶۸
- شکل ۴-۲۴: کارایی میکرو F1 روشهای فیلتر با افزایش ویژگی در مجموعه داده WEBKB..... ۶۸
- شکل ۴-۲۵: کارایی ماکرو F1 روشهای فیلتر با افزایش ویژگی در مجموعه داده WEBKB..... ۶۹
- شکل ۴-۲۶: کارایی انتخاب ویژگی تصادفی با انتخاب ۱۰۰ ویژگی..... ۷۰
- شکل ۴-۲۷: کارایی میکرو F1 انتخاب ویژگی تصادفی بر روی مجموعه داده‌های مختلف..... ۷۰
- شکل ۶-۱: کارایی روشهای فیلتر با ۱۵۰ ویژگی برتر بر روی مجموعه داده NEWSGROUP..... ۸۰
- شکل ۶-۲: دقت و بازیابی روشهای فیلتر با ۵ ویژگی برتر بر روی مجموعه داده NEWSGROUP..... ۸۰
- شکل ۶-۳: دقت و بازیابی روشهای فیلتر با ۱۵۰ ویژگی برتر بر روی مجموعه داده NEWSGROUP..... ۸۰
- شکل ۶-۴: دقت و بازیابی روشهای فیلتر با ۵۰ ویژگی برتر بر روی مجموعه داده REUTER-R8..... ۸۱
- شکل ۶-۵: دقت و بازیابی روشهای فیلتر با ۱۵۰ ویژگی برتر بر روی مجموعه داده REUTER-R8..... ۸۱
- شکل ۶-۶: دقت و بازیابی روشهای فیلتر با ۱۰۰ ویژگی برتر بر روی مجموعه داده WEBKB..... ۸۱
- شکل ۶-۷: دقت و بازیابی روشهای فیلتر با ۱۵۰ ویژگی برتر بر روی مجموعه داده WEBKB..... ۸۲

فهرست جدول‌ها

۲۸	جدول ۳-۱: مجموعه داده‌های بیوانفورماتیک
۴۷	جدول ۴-۱: مجموعه داده REUTER-R8
۴۸	جدول ۴-۲: مجموعه داده 20 NEWSGROUPS
۴۹	جدول ۴-۳: مجموعه داده WEBKB
۵۶	جدول ۴-۴: ماتریس اشتراک روش‌ها در مجموعه داده NEWSGROUP
۵۶	جدول ۴-۵: ماتریس اشتراک روش‌ها در مجموعه داده WEBKB
۵۷	جدول ۴-۶: ماتریس اشتراک روش‌ها در مجموعه داده REUTER-R8
۸۲	جدول ۶-۱: ماتریس اشتراک روش‌ها در ۱۰ ویژگی در مجموعه داده NEWSGROUP
۸۲	جدول ۶-۲: ماتریس اشتراک روش‌ها در ۵۰ ویژگی در مجموعه داده NEWSGROUP
۸۲	جدول ۶-۳: ماتریس اشتراک روش‌ها در ۱۰ ویژگی در مجموعه داده REUTER
۸۳	جدول ۶-۴: ماتریس اشتراک روش‌ها در ۵۰ ویژگی در مجموعه داده REUTER
۸۳	جدول ۶-۵: ماتریس اشتراک روش‌ها در ۱۰ ویژگی در مجموعه داده WEBKB
۸۳	جدول ۶-۶: ماتریس اشتراک روش‌ها در ۵۰ ویژگی در مجموعه داده WEBKB

۱-۱ - سابقه و انگیزه تحقیق

امروزه، وجود دستگاه‌ها، وسایل و حسگرهای قوی در کاربردهای مختلف برای نمونه‌برداری و استخراج مشخصات و ویژگی‌های نمونه‌های داده‌ای باعث به وجود آمدن مجموعه‌هایی با تعداد ویژگی‌های بسیار زیاد^۱ شده است. از طرفی سیستم‌های کاربردی مختلف در دنیای واقعی نیازمند سرعت و دقت بالایی می‌باشند که با این حجم عظیم داده‌ها کارایی خود را از دست می‌دهند. از اینرو یادگیری ساختار داده و انتخاب ویژگی‌ها و مشخصاتی که نمایانگر این ساختار باشند بسیار حایز اهمیت بوده و خصوصاً در دهه اخیر به یکی از حوزه‌های فعال تحقیقاتی در نقاط مختلف دنیا تبدیل شده است. مثلاً با رشد روزافزون تکنولوژی اینترنت شاهد به وجود آمدن مجموعه داده‌های بسیار بزرگ از سندهای الکترونیکی هستیم. طبقه‌بندی این سندهای الکترونیکی که شامل هزاره‌ها و هزاره‌ها به زبان‌های مختلف می‌باشند از چالش‌های مهم سالهای اخیر می‌باشد که در این راستا روش‌های زیادی برای انتخاب واژه‌های مناسب به منظور کاهش فضای ویژگی و طبقه‌بندی مناسب‌تر انجام گرفته است [1]. انتخاب ویژگی در حوزه‌های پزشکی که شامل انتخاب ژن^۲ های موثر در تشخیص بیماری‌ها می‌باشند نمونه‌ای دیگر از کاربردهای جدید انتخاب ویژگی در داده‌های ابعاد بالا در دهه اخیر می‌باشند [3], [2]. در سرویس‌های وب، انتخاب ویژگی برای شناسایی نیاز کاربران، تشخیص هرزنامه‌ها، بهبود جستجو و غیره مورد استفاده قرار می‌گیرد (مانند، [5], [4]). لذا حجم بالای تحقیقات انجام شده در سالهای اخیر بر روی انتخاب ویژگی در حوزه‌های مختلف نشان‌دهنده اهمیت هرچه بیشتر این موضوع می‌باشد.

1-High Dimensional Data

² Gene

برای انتخاب ویژگی در حوزه‌های مختلف خصوصا در زمینه اسناد الکترونیکی و داده‌های بیوانفورماتیک که شامل مجموعه‌هایی با ابعاد بسیار بالا می‌باشند راهکارهای مختلفی ارائه شده است. که این راهکارها بر اساس ماهیت داده‌های مورد بررسی تعریف می‌شوند. در مجموعه داده‌های متنی، ویژگی‌ها معمولا خصوصیات آماری می‌باشند که از واژه‌های موجود در اسناد بدست می‌آیند. بنابراین ویژگی‌ها در این حوزه ماهیتی گسسته دارند. روش‌های بهره اطلاعات^۱ [6]، شاخص بهبودیافته جینی^۲ [7] و اطلاعات متقابل^۳ [8] از روش‌های پرکاربرد انتخاب ویژگی در این حوزه می‌باشند. بر خلاف مجموعه داده‌های متنی، در مجموعه داده‌های میکروآرایه (بیوانفورماتیک) ویژگی‌ها که مقدار بیان ژنها می‌باشند مقداری حقیقی داشته به طوری که فضای پیوسته ای را به وجود می‌آورند. الگوریتم‌های انتخاب ویژگی مختلفی در این حوزه ارائه شده است که برخی از آنها مانند آزمایش خی دو^۴ [6] و الگوریتم رتبه‌دهی لاپلاسی^۵ [۹] می‌باشند. در این پایان‌نامه الگوریتم‌هایی برای انتخاب ویژگی هر دو حوزه ارائه می‌شود.

انتخاب ویژگی در حوزه‌های مختلف می‌تواند با استخراج ویژگی نیز همراه باشد. استخراج ویژگی، تولید یک فضای ویژگی جدید بر اساس ویژگی‌های اصلی می‌باشد [۱۰]. یکی از روش‌های استخراج ویژگی روش تحلیل مولفه‌های سازنده^۶ [12], [11] می‌باشد. به طور کلی، انتخاب و استخراج ویژگی و یا کاهش ابعاد داده‌ای (کاهش تعداد ویژگی‌ها) می‌تواند به عنوان پیش پردازش برای طبقه‌بندی، خوشه‌بندی و بازیابی اطلاعات، به منظور سرعت بخشیدن و افزایش دقت این عملیات مورد استفاده قرار گیرد. روش‌های انتخاب ویژگی از یک دیدگاه کلی قابل دسته‌بندی به روش‌های باناظر^۷ و بی‌ناظر^۸ می‌باشند. در روش‌های باناظر معمولا اهمیت ویژگی با توجه به همبستگی^۹ که ویژگی به برجسب کلاس^{۱۰} خود دارد، مورد ارزیابی قرار می‌گیرد. برخی از روش‌های کلاسیک انتخاب ویژگی، روش Lasso توسط [۴]، روش Relief توسط [۱۳] می‌باشد. در روش Lasso با استفاده از یک تابع

¹ Information Gain

² GINI Index

³ Mutual Information

⁴ Chi-Square Test

⁵ Laplacion Score

⁶ Principal Component Analysis

⁷ Supervised

⁸ Unsupervised

⁹ correlation

¹⁰ Class label

بهینه سازی سعی در انتخاب ویژگی‌ها می‌شود. در این روش از ضرایب رگرسیون برای تولید یک راه‌حل تنک استفاده می‌نماید. در نهایت ویژگی‌هایی که ضریب رگرسیون صفر دارند حذف می‌شوند. دسته دیگر از الگوریتم‌های انتخاب ویژگی، روش‌های بی‌ناظر، شامل روش‌هایی می‌شوند که در آنها "هدف" انتخاب ویژگی‌هایی است که زیرساخت هندسی¹ موجود در ساختار داده را به‌بهترین وجه حفظ نمایند. یکی از ساده‌ترین این روش‌ها، روش واریانس داده‌ای² می‌باشد. در این روش، داده‌ها بر روی بعدی که بیشترین واریانس را دارد نگاشت³ شده و انتخاب ویژگی از ویژگی‌های اصلی انجام می‌گیرد. روش تحلیل مولفه‌های اصلی توسط [11] نیز از حداکثر واریانس استفاده می‌نماید اما در ادامه ویژگی‌ها را تبدیل و انتخاب ویژگی را از ویژگی‌های تبدیل‌یافته⁴ انجام می‌دهد. اخیراً الگوریتم‌های انتخاب ویژگی⁵ LapAOFS با بهره‌گیری از مفاهیم طراحی تجربی توسط [14] ارائه شده است. این الگوریتم با حداقل سازی ماتریس کوواریانس داده‌ای سعی در بدست آوردن ویژگی‌هایی نموده اند که هم ساختار هندسی و هم ساختار فضایی⁶ داده‌ها را حفظ نمایند. در این الگوریتم که برای انتخاب ویژگی در داده‌های ابعاد بالا می‌باشد امکان استفاده از برچسب کلاس‌ها در صورت موجود بودن وجود ندارد.

۱-۲- مسئله‌ی تحقیق و اهداف پایان‌نامه

به طور کلی در این پایان‌نامه دو روش انتخاب ویژگی مجزا برای داده‌های ابعاد بالای میکروآرایه و متن‌های الکترونیکی ارائه می‌شود.

روش انتخاب ویژگی LapAOFS که از آخرین روش‌های بی‌ناظر در انتخاب ویژگی می‌باشد. در این روش همانطور که در مقدمه بدان اشاره شد امکان استفاده از برچسب کلاس‌ها در صورت موجود بودن، وجود ندارد. استفاده از برچسب کلاس‌ها و ساختار داده‌ای به طور همزمان می‌تواند به مراتب الگوریتم‌های بی‌ناظر را در انتخاب ویژگی کارا تر سازد. چرا که استفاده از اطلاعات خبره علاوه بر خصوصیات نهفته در ساختار داده می‌تواند باعث هدایت انتخاب ویژگی گردد [15]. لذا در این پایان‌نامه با تکیه بر یک الگوریتم استخراج ویژگی سعی می‌شود که فضای داده‌ای جدیدی را برای

¹ underlying geometrical structure

² data variance

³ project

⁴ transformed

⁵ Laplacian A-Optimality Feature Selection

⁶ manifold

الگوریتم رتبه‌دهی بی‌ناظر LapAOFS فراهم می‌آوردیم. در روش پیشنهادی با کاهش تاثیر نقاط داده‌ای و با در نظر گرفتن مرزهای ویژگی در هر کلاس سعی در تاثیر دادن برچسب کلاس‌ها در انتخاب ویژگی این الگوریتم شده است.

روش دوم در این پایان‌نامه به بررسی انتخاب ویژگی و روش‌های احتمالاتی انتخاب ویژگی مبتنی بر فیلتر در مجموعه داده‌های متنی می‌پردازد. در داده‌های ابعاد بالای متن مساله انتخاب واژه‌های موثر در طبقه‌بندی اسناد با روش‌های احتمالاتی مناسب از چالش‌های مهم در ارزیابی روش‌های انتخاب ویژگی مبتنی بر فیلتر می‌باشد. روش‌های احتمالاتی مختلفی در این حوزه مطرح شده‌اند که از عامل‌های احتمالاتی متفاوتی استفاده می‌کنند. برای اولین بار در این پایان‌نامه عامل‌های احتمالاتی^۱ روش‌های انتخاب ویژگی احتمالاتی مبتنی بر فیلتر مورد بررسی قرار می‌گیرند به طوریکه در نهایت با بکارگیری مهمترین عامل‌های تاثیرگذار در انتخاب ویژگی‌های مناسب، در میان روش‌های احتمالاتی مختلف، به ارزیابی روشی احتمالاتی پرداخته می‌شود.

۳-۱- دورنمای پایان‌نامه

این پایان‌نامه با احتساب این فصل در ۵ فصل کلی گردآوری شده است. در ذیل به طور خلاصه به تشریح هر فصل می‌پردازیم:

فصل اول به اهمیت موضوع، چالش‌ها و اهداف پایان‌نامه در کنار معرفی برخی از آخرین روش‌های ارزیابی شده در حوزه مورد بررسی می‌پردازد،

فصل دوم به مروری بر کارهای انجام شده در زمینه انتخاب ویژگی می‌پردازد. در این فصل تمامی روش‌های انتخاب ویژگی که در فصل‌های بعد از آنها استفاده شده، تشریح می‌شود،

فصل سوم به ارزیابی روش انتخاب ویژگی بی‌ناظر برای ویژگی‌هایی که مقادیر حقیقی دارند می‌پردازد. در این فصل با تکیه بر یک الگوریتم استخراج ویژگی پیشنهادی سعی در تاثیر دادن برچسب کلاس در الگوریتم انتخاب ویژگی LapAOFS می‌کند. سپس روشی باناظر به منظور بررسی بیشتر استخراج ویژگی ارزیابی شده، مورد آزمایش قرار می‌گیرد. از مجموعه داده‌های ابعاد بالای بیوانفورماتیک برای ارزیابی روش پیشنهادی استفاده می‌گردد،

^۱ Probabilistic elements

فصل چهارم به ارایه روش انتخاب ویژگی احتمالاتی مبتنی بر فیلتر برای مجموعه داده‌های متنی می‌پردازد. عامل‌های احتمالاتی روش‌های فیلتر مختلف مورد بررسی قرار گرفته و بر اساس تاثیرگذارترین عامل‌ها روش فیلتر احتمالاتی جدیدی ارایه می‌شود. سپس آزمایشات، مقایسه‌ها و تحلیل‌های بسیاری بر روی سه مجموعه داده متنی مشهور انجام می‌گیرد.

فصل پنجم در نهایت به جمع‌بندی و کارهای تحقیقاتی بیشتر در آینده می‌پردازد.

۲-۱- مقدمه

انتخاب ویژگی فرآیندی است که در آن زیرمجموعه‌ای از ویژگی‌های موجود بر اساس معیارهای اهمیت ویژگی انتخاب می‌شوند [۱۶]. بعد از اولین مقالات در زمینه انتخاب ویژگی‌های مرتبط توسط [۱۷]، [۱۸] تا به امروز پیشرفت‌های زیادی در این امر حاصل شده است. این پیشرفت‌ها از یکسو حاصل تلاش محققان در زمینه‌های داده‌کاوی و از سوی دیگر حاصل پیشرفت تکنولوژی داده‌کاوی در زمینه‌های علمی مختلف می‌باشد. انتخاب ویژگی‌ها در زمینه‌های مختلف غالباً به منظور پیش‌بینی، طبقه‌بندی، خوشه‌بندی و کاهش حجم داده‌ها انجام می‌گیرد.

از سوی دیگر انتخاب ویژگی با توجه به ماهیت ویژگی مانند ترتیبی، گسسته یا پیوسته بودن، با روش‌های مختلفی انجام می‌گیرد. مثلاً در داده‌های بیوانفورماتیک ویژگی‌ها مقداری حقیقی داشته و از این‌رو حوزه انتخاب ویژگی به صورت پیوسته مورد بحث قرار می‌گیرد. در عین حال انتخاب ویژگی در زمینه متن‌های اینترنتی به صورت انتخاب ویژگی بر اساس اطلاعات آماری از ویژگی‌ها مانند تعداد واژه در کلاس بدست می‌آیند که مشخصاً فضای انتخابی گسسته را به وجود می‌آورند.

۲-۲- روش‌های انتخاب ویژگی

طبقه‌بندی داده‌های ابعاد بالا در مجموعه داده‌های مختلف خود شامل مسائل و مشکلات زیادی است که باعث ارائه‌ی روش‌های مختلف انتخاب ویژگی در حوزه‌های متعدد شده است. انتخاب ویژگی به طور کلی قابل دسته‌بندی به دو روش باناظر و بی‌ناظر می‌باشد. در روش‌های باناظر اهمیت ویژگی‌ها با توجه به همبستگی آنها با برجسب کلاس‌ها مورد ارزیابی قرار می‌گیرد. از روش‌های متداول انتخاب ویژگی باناظر می‌توان به روش‌های ReliefF توسط [۱۳]، Lasso توسط [۱۹]، SVM-RFE توسط [۲۰] و Fisher score در کتاب طبقه‌بندی الگو توسط [۲۱] اشاره داشت.

در عین حال استفاده از روش‌های باناظر برای بسیاری از مجموعه داده‌ها که بدون برچسب هستند ممکن نیست. از این رو محققان به ارائه‌ی روش‌های یادگیری بی‌ناظر برای طبقه‌بندی داده‌ها و انتخاب ویژگی نموده‌اند. توجه به این نکته حائز اهمیت است که روش انتخاب ویژگی بی‌ناظر روشی است که به برچسب کلاس توجهی ندارد، لذا این روش‌ها را می‌توان برای هر نوع مجموعه داده‌ای بکار برد. برخی از این روش‌ها عبارتند از:

روش انتخاب ویژگی بی‌ناظر ترتیبی [۱۶]، روش Laplacian score توسط [۲۲]، روش انتخاب ترکیبی بر اساس حداکثرسازی احتمال [۲۳]، روش انتخاب حداکثر آنتروپی [۲۴] و روش تحلیل مؤلفه‌های سازنده (PCA) [۲۵]. در روش تحلیل مؤلفه‌های اصلی ابتدا عمل تبدیل ویژگی‌ها به فضای دیگر انجام گرفته و در نهایت مجموعه‌ای از ویژگی‌های تبدیل یافته به جای ویژگی‌های اصلی معرفی می‌شوند. اخیراً روش LapAOFs توسط [۱۴] بر اساس رتبه‌دهی لاپلاسی [۲۶] و نسخه ارتقا یافته آن [۲۲] ارائه شده است. این روش که راه‌کاری بی‌ناظر برای رتبه‌دهی ویژگی‌ها می‌باشد، بر اساس ساختار داده‌ها در گراف سعی در انتخاب موثرترین ویژگی‌ها می‌نماید.

۲-۳- انتخاب ویژگی در متن

یکی از زمینه‌های بسیار پرکاربرد انتخاب ویژگی در حوزه طبقه‌بندی متن‌ها^۱ می‌باشد. طبقه‌بندی متن فرآیند طبقه‌بندی یک سند بر اساس محتویات درونی آن می‌باشد [۲۷]. در واقع با پیشرفت روزافزون تکنولوژی اینترنت حجم سندهای الکترونیکی به طور چشم‌گیری در حال افزایش می‌باشد. لذا، طبقه‌بندی این سندها بر اساس ویژگی‌های منتخب به جای استفاده از حجم انبوهی از ویژگی‌ها از مهم‌ترین مسائل داده‌کاوی اخیر می‌باشد. در دهه‌ی اخیر، روش‌های زیادی برای انتخاب و طبقه‌بندی این داده‌ها توسط محققان ارائه شده است. روش‌های انتخاب ویژگی قابل دسته‌بندی به روش‌های فیلتر^۲، روپوش^۳ و ترکیبی^۴ و ادغام شده^۵ می‌باشند [۲۸]. در بخش‌های زیر تشریحی مختصر از هر یک از این روش‌ها آورده شده است.

^۱ Text Classification

^۲ Filter

^۳ Wrapper

^۴ Hybrid

^۵ Embedded

۲-۳-۱- روش روپوش

در روش روپوش، یک ماشین یادگیر دلخواه برای رتبه‌دهی به زیرمجموعه‌ای از ویژگی‌ها، با توجه به توان آن ویژگی‌ها در امر پیش‌بینی به کار می‌رود. در این روش معمولاً روش‌های اکتشافی^۱ مانند جستجوی تابو^۲، الگوریتم ژنتیک^۳، کلونی مورچه‌ها^۴ و غیره، همراه با روش‌های فیلتر مورد استفاده قرار می‌گیرند. این روش‌ها هنگامی که فضای ویژگی بسیار بزرگ باشد با جستجوی مبتنی بر طبیعت قادر به یافتن راه حل‌های بهینه و یا نزدیک به بهینه می‌باشند.

اخیراً و با افزایش حجم داده‌ها در مجموعه داده‌های بزرگ از روش‌های روپوش استفاده زیادی شده است که حجم عظیم مقالات در این زمینه گواهی این مدعاست. در [۲۹] راهکاری کلی بر مبنای روش روپوش برای انتخاب ویژگی‌هایی که بیشترین وابستگی به کلاس داده را دارند، معرفی شده است. در [۲۳]، [۳۰]، [۳۱] روش‌هایی بر اساس الگوریتم ژنتیک برای دسته‌بندی متن ارائه شده است. نسخه تغییر یافته‌ای بر مبنای الگوریتم کلونی مورچه‌ها در [۳۲] استفاده شده است که سعی در حداکثر استفاده از فضای جستجوی ویژگی‌ها برای یافتن ویژگی‌های موثر نموده است.

۲-۳-۲- انتخاب ویژگی مبتنی بر فیلتر

انتخاب ویژگی مبتنی بر فیلتر نیازمند تحلیل آماری مجموعه ویژگی‌ها بدون استفاده از هر نوع مدل و یا ماشین یادگیر می‌باشد [۳۳]. روش‌های فیلتر نسبت به بقیه‌ی روش‌ها روش‌های سریع‌تری هستند اما از طرف دیگر وابستگی بین ویژگی‌ها را مورد بررسی قرار نمی‌دهند. برخی از موفق‌ترین روش‌های مبتنی بر فیلتر در حوزه متن عبارتند از:

توان واژه^۵ [۳۴]، روش احتمالاتی DFS [۳۵]، اطلاعات متقابل [31]، آزمایش خی‌دو [۸]، بهره‌ی اطلاعاتی [۶]، اندیس بهبود یافته‌ی جینی [۷]، معیار انحراف از توزیع پواسون [۳۶]، حداقل اختلاف کلاس^۶ [۳۷]، میزان ابهام^۷ [۳۸]، میزان تمایز دهندگی بین کلاس‌ها^۸ [۳۹] و آزمون فرض

¹ Heuristics

² Tabu Search

³ Genetic Algorithm

⁴ Ant Colony

⁵ Term Strength

⁶ minimum class difference

⁷ ambiguity measure

⁸ class discriminating measure