

تأثیر ترکیب روش‌های انتخاب ویژگی و بسته‌بندی در بهبود پیش‌بینی اشکال نرم‌افزار

فاطمه علیقارداشی^۱، دانشجوی کارشناسی ارشد؛ محمدعلی زارع چاهوکی^۲، استادیار

۱- دانشکده مهندسی برق و کامپیوتر - دانشگاه یزد - یزد - ایران - f.alighardashi@stu.yazd.ac.ir

۲- دانشکده مهندسی برق و کامپیوتر - دانشگاه یزد - یزد - ایران - chahooki@yazd.ac.ir

چکیده: حفظ کیفیت محصول نرم‌افزاری با آزمون‌های دوره‌ای قبل از نصب، یکی از پرهزینه‌ترین فعالیت‌ها در پروژه‌های فناوری اطلاعات است. با توجه به منابع محدود برای آزمون ماژول‌ها در پروژه‌های نرم‌افزاری، بهتر است ابتدا ماژول‌های مستعد اشکال شناسایی شوند و منابع آزمون در جهت شناسایی اشکال در این ماژول‌ها متمرکز گردند. پیش‌بینی‌کننده‌های اشکال مبتنی بر الگوریتم‌های یادگیری ماشین، ابزارهای مقرون‌به‌صرفه‌ای برای شناسایی ماژول‌های مستعد اشکال هستند. پژوهش‌های گسترده‌ای در این حوزه برای یافتن ارتباط بین ویژگی‌های ماژول‌های نرم‌افزاری و مستعد اشکال بودن آن‌ها صورت پذیرفته است. برخی از این ویژگی‌ها در الگوریتم‌های پیش‌بینی‌کننده به‌گونه‌ای هستند که نه تنها سبب بهبود دقت در فرآیند یادگیری نمی‌شوند بلکه کاهش دقت را نیز در پی خواهند داشت. در این پژوهش با توجه به عملکرد خوب روش انتخاب ویژگی روبه‌جلو در انتخاب ویژگی‌های مؤثر، زیرمجموعه اولیه در این روش با استفاده از تلفیق ویژگی‌های با رتبه بالا در روش‌های مختلف فیلتر انتخاب می‌شود. روش پیشنهادی علاوه بر بهبود دقت سبب افزایش سرعت همگرایی در انتخاب ویژگی می‌شود. نتایج حاصل از پیاده‌سازی و ارزیابی نتایج تجربی به‌دست‌آمده در دادگان ناسا با معیار AUC، بیانگر مؤثر بودن روش پیشنهادی در بهبود دقت و سرعت پیش‌بینی ماژول‌های نرم‌افزاری مستعد اشکال است.

واژه‌های کلیدی: پیش‌بینی اشکال نرم‌افزار، یادگیری ماشین، انتخاب ویژگی، فیلتر، بسته‌بندی.

The Effectiveness of the Combination of Filter and Wrapper Feature Selection Methods to Improve Software Fault Prediction

F. Alighardashi¹, MSc Student; M. A. Zare Chahooki², Assistant Professor

1- Faculty of Electrical and Computer Engineering, University of Yazd, Yazd, Iran, Email: f.alighardashi@stu.yazd.ac.ir

2- Faculty of Electrical and Computer Engineering, University of Yazd, Yazd, Iran, Email: chahooki@yazd.ac.ir

Abstract: Improving the software product quality before releasing by periodic tests is one of the most expensive activities in software projects. Due to limited resources to test modules in software projects, it is important to identify fault-prone modules and use the test sources for fault prediction in these modules. Software fault predictors based on machine learning algorithms, are effective tools for identifying fault-prone modules. Extensive studies are being done in this field to find the connection between features of software modules, and fault-prone have been done. Some of these features in predictive algorithms are like that not only did not improve the accuracy of the learning process, but also will be reduced the accuracy. In this study, due to the excellent performance of Forward Feature Selection (FS) method for effective selection of features, the initial subset of this method has been selected by using of combination of high ranking features in different Filter methods. The proposed method causes increment the speed of the convergence of feature selection as well as the accuracy improvement. The obtained results on NASA dataset with AUC criteria, indicates the effectiveness of this method in the improvement of the accuracy and the speed of software fault prediction.

Keywords: Software fault prediction, machine learning, feature selection, filter, wrapper.

تاریخ ارسال مقاله: ۱۳۹۴/۱۱/۱۸

تاریخ اصلاح مقاله: ۱۳۹۵/۰۱/۲۹ و ۱۳۹۵/۰۲/۲۲

تاریخ پذیرش مقاله: ۱۳۹۵/۰۳/۱۱

نام نویسنده مسئول: محمدعلی زارع چاهوکی

نشانی نویسنده مسئول: ایران - یزد - بلوار دانشگاه - دانشگاه یزد - دانشکده مهندسی برق و کامپیوتر.

۱- مقدمه

مدل‌های پیش‌بینی اشکال نرم‌افزار با شروع از سال ۱۹۹۰ تا به امروز توسط پژوهشگران بسیاری مورد مطالعه قرار گرفته‌اند که نتیجه آن شناسایی ماژول‌های نرم‌افزاری مستعد اشکال قبل از آزمون سیستم با استفاده از این مدل‌ها بوده است. در پژوهش [۵]، منزیس و همکارانش میزان تشخیص اشکال در ماژول‌های نرم‌افزاری توسط مدل‌های پیش‌بینی را حدود ۷۱ درصد اعلام کردند که بالاتر از میزان تشخیص اشکال توسط منابع انسانی (۶۰ درصد) بوده است. ایشان در جلساتی که کنفرانس متریک IEEE در سال ۲۰۰۲ برگزار کرد، ادعای فاگان مبنی بر این که بررسی کد نرم‌افزار قبل از فرآیند آزمون، توسط منابع انسانی می‌تواند ۹۵ درصد از نقایص نرم‌افزار را برطرف کند را رد کرد. چراکه بر اساس پژوهش مذکور ثابت شد که نسبت تشخیص اشکال‌های کد نرم‌افزاری توسط منابع انسانی حدود ۶۰ درصد می‌شود. منزیس در این راستا بیان کرد که هر یک از اعضای تیم بررسی کد نرم‌افزار می‌توانند حدود ۸ تا ۲۰ خط کد را در هر دقیقه بررسی کنند و آن‌ها فرآیند بررسی را به‌وسیله تمام اعضای تیم متشکل از ۴ تا ۶ نفر انجام داده‌اند. اما بررسی کد بسیار زمان‌بر بوده و تنها حدود ۶۰ درصد از اشکال‌های موجود در ماژول‌های نرم‌افزاری تشخیص داده شده‌اند. بدین سبب، ثابت شد که استفاده از مدل‌های پیش‌بینی اشکال نرم‌افزار در مقایسه با بررسی نرم‌افزار توسط تیم نرم‌افزاری، بسیار مقرون‌به‌صرفه‌تر است [۴].

انتخاب ویژگی، یکی از رویکردهای بهبود دقت و سرعت در روش‌های یادگیری ماشین است [۶]. در سال‌های اخیر پژوهش‌های متعددی در حوزه انتخاب ویژگی در پیش‌بینی اشکال نرم‌افزار صورت پذیرفته است [۷]. پژوهش‌های مذکور نشان‌دهنده تعداد بالای ویژگی‌های استخراج‌شده از ماژول‌های نرم‌افزاری هستند. نتایج پژوهش‌ها در حوزه کاهش ویژگی‌ها نشان داده است که انتخاب مجموعه‌ای کمتر از ویژگی‌های اولیه می‌تواند سبب افزایش دقت در الگوریتم‌های یادگیری ماشین شود. پیش‌بینی اشکال نرم‌افزار با استفاده از لیست کاملی از ویژگی‌های در دسترس، اغلب سبب کاهش دقت در پیش‌بینی می‌شود. چراکه در برخی موارد ممکن است که برخی از این ویژگی‌ها نامربوط یا افزونه باشند و سبب ایجاد سردرگمی مدل یادگیری و نیز پیچیدگی‌های اضافی در فرآیند ساخت مدل پیش‌بینی توسط روش یادگیری ماشین شوند. برای مقابله با این مسئله، از انتخاب ویژگی استفاده می‌شود. انتخاب ویژگی فرایند یافتن زیرمجموعه‌ای از ویژگی‌هاست که با آن‌ها یک مدل، دقتی مشابه یا بهتر نسبت به استفاده از تمام ویژگی‌ها دارد [۷].

یکی از مشکلات دادگان‌های ارزیابی این است که در بیش‌تر مواقع تمام ویژگی‌های دادگان برای یافتن دانشی که در آن‌ها نهفته است مهم و حیاتی نیستند، این مسئله در پژوهش‌های بسیاری [۸-۱۰] به‌عنوان فلسفه اصلی نیاز به کاهش ویژگی‌های دادگان برای ایجاد مدل یادگیری بهتر و بهبود دقت مدل بر دادگان بیان شده است. چراکه روش‌های یادگیری ماشین، با استفاده از ارتباط‌های بین ویژگی‌های

در طول سال‌های اخیر درخواست برای کیفیت نرم‌افزار رشد بالایی داشته است. همچنین با توجه به پیچیده‌تر شدن سیستم‌های نرم‌افزاری، پیچیدگی مسائل مرتبط به آزمون نرم‌افزار نیز به‌طور فزاینده‌ای افزایش یافته است. قابلیت اندازه‌گیری مستعد اشکال بودن ماژول‌های نرم‌افزاری می‌تواند به‌شدت روی کاهش هزینه و بهبود فرآیند آزمون نرم‌افزار تأثیرگذار باشد. اکثر خطاهای سیستم نرم‌افزاری از تعداد اندکی از اجزای سیستمی ناشی می‌شوند. قانون ۸۰:۲۰ در این زمینه بیانگر آن است که ۲۰ درصد از ماژول‌های نرم‌افزاری سبب ایجاد ۸۰ درصد خطا، هزینه و کار مجدد می‌شود. بنابراین با تلاش برای رفع خطا و بهبود کیفیت بخش کوچکی از نرم‌افزار، اتلاف زمان کمتری در کل پروژه تولید خواهیم داشت [۱].

با توجه به منابع محدود برای آزمون ماژول‌ها در پروژه‌های نرم‌افزاری، بهتر است ابتدا ماژول‌های مستعد اشکال (fp^1) شناسایی شده و از ماژول‌های غیرمستعد اشکال (nfp^2) تفکیک شوند و منابع آزمون در جهت شناسایی اشکال در ماژول‌های مستعد اشکال متمرکز شود. این مسئله موجب مطرح شدن مبحث پیش‌بینی ماژول‌های نرم‌افزاری مستعد اشکال شد که به‌اصطلاح «پیش‌بینی اشکال نرم‌افزار» نیز گفته می‌شود [۱].

پیش‌بینی اشکال نرم‌افزار، مدیران پروژه‌های نرم‌افزاری را قادر می‌سازد که تلاش برای بهبود کیفیت نرم‌افزار را در قسمت‌های موردنیاز متمرکز کنند. به‌عنوان مثال، قبل از آزمون سیستم، شناسایی اجزایی که به احتمال زیاد در حین عملیات باعث بروز خطا می‌گردند، می‌تواند اثربخشی تلاش برای آزمون نرم‌افزار را بهبود بخشد. به همین دلیل روش‌های مختلف مدل‌سازی پیش‌بینی اشکال نرم‌افزار توسعه داده شده‌اند [۲].

پیش‌بینی‌کننده‌های اشکال مبتنی بر الگوریتم‌های یادگیری ماشین، ابزارهای مقرون‌به‌صرفه‌ای برای شناسایی ماژول‌های نرم‌افزاری مستعد اشکال هستند. پیش‌بینی اشکال نرم‌افزار مبتنی بر یادگیری ماشین رویکردی است که از ویژگی‌های نرم‌افزاری پروژه‌های پیشین استفاده کرده سپس مدلی ایجاد می‌کند که ماژول‌های مستعد اشکال نسخه بعدی نرم‌افزار را پیش‌بینی می‌کند [۳].

مدل مستعد اشکال بر اساس اطلاعات کدهای نرم‌افزاری و وجود یا عدم وجود اشکال در آن‌ها ساخته شده است. اطلاعات اشکال به‌عنوان متغیر وابسته در مدل‌های پیش‌بینی اشکال نرم‌افزار در نظر گرفته می‌شود. اگر خطایی در حین انجام آزمایش سیستم نرم‌افزاری گزارش شود، اطلاعات اشکال ماژول با مقدار غیرصفر مقدرده می‌شود و در غیراین‌صورت مقدار صفر یعنی بدون اشکال به آن نسبت داده می‌شود. برای مدل‌سازی پیش‌بینی، ویژگی‌های نرم‌افزار به‌عنوان متغیر مستقل در مدل پیش‌بینی استفاده شده و اطلاعات اشکال به‌عنوان متغیر وابسته یا هدف در مدل پیش‌بینی استفاده می‌شود [۴].

در برخی دیگر دقت بسیار پایینی داشته‌اند. بنابراین نمی‌توان تعیین کرد که کدام روش فیلتر برای پیش‌بینی اشکال پروژه نرم‌افزاری جدید بهتر است. همان‌گونه که در نتایج آزمایش‌های تجربی نشان داده شده است، روش تلفیقی ارائه‌شده در این پژوهش در تمامی دادگان‌های آزمون شده مؤثر است. به عبارت دیگر نتایج به‌دست‌آمده از پیاده‌سازی روش پیشنهادی نشان‌دهنده بهبود دقت و سرعت نسبت به استفاده از تمامی ویژگی‌ها است.

همچنین هرچند تاکنون از SVM [۲] در پژوهش‌های مختلفی در حوزه پیش‌بینی اشکال نرم‌افزار استفاده شده است ولی از برخی از توسعه‌های جدیدی که در این فرضیه یادگیری داده شده است، در حوزه پیش‌بینی اشکال نرم‌افزار استفاده نشده است. نوآوری جنبی ارائه‌شده در این پژوهش، افزایش دقت دسته‌بندی با ماشین بردار پشتیبان دوقلو هموارشده با حاشیه پارامتریک (STPMSVM) [۱۷] است.

در ادامه این پژوهش در بخش ۲، روش‌های انتخاب ویژگی معرفی شده‌اند. در بخش ۳، پژوهش‌های مبتنی بر انتخاب ویژگی در پیش‌بینی اشکال نرم‌افزار مرور شده‌اند. در بخش ۴، روش پیشنهادی این پژوهش بیان شده است. در بخش ۵ و ۶ به ترتیب دادگان و معیارهای ارزیابی کارایی در حوزه پیش‌بینی اشکال مبتنی بر یادگیری ماشین آورده شده است. در بخش ۷ نیز نتایج تجربی آزمایش‌های صورت‌پذیرفته جهت بررسی تأثیر روش پیشنهادی بر پیش‌بینی ماژول‌های مستعد اشکال، آورده شده است و این نتایج با نتایج دیگر پژوهش‌های انجام‌شده در زمینه پیش‌بینی اشکال نرم‌افزار مقایسه و تحلیل شده است. در نهایت در بخش ۸ نیز نتیجه‌گیری و پژوهش‌های پیش رو آورده شده است.

۲- انتخاب ویژگی

کاهش ویژگی، یکی از پرکاربردترین روش‌هایی است که در مبحث یادگیری ماشین مطرح شده است. هدف کاهش ویژگی، حذف ویژگی‌های غیرضروری و انتخاب مناسب‌ترین ویژگی‌ها از بین مجموعه ویژگی‌های اولیه، برای افزایش عملکرد الگوریتم‌های یادگیری است [۱۸]. معمولاً جستجوی جامع برای پیدا کردن مناسب‌ترین ویژگی‌ها به‌لحاظ هزینه محاسباتی غیرممکن است. بنابراین کاهش ویژگی تبدیل به یک چالش عمده در شناسایی الگو و یادگیری ماشین شده است. این روش در بسیاری از کاربردها از جمله طبقه‌بندی و رگرسیون اهمیت بسیاری دارد. چراکه در این کاربردها معمولاً تعداد زیادی ویژگی وجود دارد، که بسیاری از آن‌ها یا اثربخش نیستند و یا سبب کاهش دقت یادگیری می‌شوند. حذف کردن این ویژگی‌ها علاوه‌بر افزایش دقت، سبب کاهش پیچیدگی محاسباتی می‌شود [۱۹].

روش‌های انتخاب ویژگی، روش‌هایی مبتنی بر کاهش ویژگی هستند. این روش‌ها سعی می‌کنند با انتخاب زیرمجموعه‌ای از ویژگی‌های اولیه، ابعاد داده‌ها را کاهش دهند. در این روش‌ها به‌دنبال

تشریح‌کننده نمونه‌های دادگان آموزشی، الگویی را استخراج می‌کنند و با استفاده از این الگو و ویژگی‌های تشریح‌کننده نمونه جدید، پیش‌بینی اشکال نرم‌افزار را برای این نمونه انجام می‌دهند [۷]. بنابراین، بی‌شک ویژگی‌های تشریح‌گر نمونه‌ها تأثیر فوق‌العاده‌ای روی دقت روش‌های یادگیری ماشین دارند و دقت و درستی مدل مورداستفاده، به ویژگی‌های مجموعه دادگان بستگی دارد. بنابراین نیاز است که بهترین ویژگی‌ها برای ایجاد مدل بهتر توسط روش‌های انتخاب ویژگی انتخاب شوند.

الگوریتم انتخاب ویژگی می‌تواند تلفیقی از تکنیک‌های جستجو برای ارائه زیرمجموعه‌ای جدید از ویژگی‌ها، همراه با معیاری برای ارزیابی که زیرمجموعه‌های متمایز را امتیازدهی می‌کند، باشد [۱۱]. رویکردی که اخیراً در حوزه انتخاب ویژگی مطرح شده است، ترکیب روش‌های انتخاب ویژگی است. ایده ترکیب روش‌های انتخاب ویژگی برای کاربردهای دیگر [۱۴-۱۲] نیز مطرح شده است. نتایج پژوهش‌های اخیر بیانگر تأثیر مثبت روش‌های ترکیبی انتخاب ویژگی در کاربردهای حوزه مهندسی نرم‌افزار از جمله پیش‌بینی اشکال نرم‌افزار نیز بوده است [۱۵، ۱۶]. بنابراین فرضیه‌ای که در این پژوهش مطرح شده است، آن است که استفاده از این روش‌ها در تشخیص ماژول‌های نرم‌افزاری مستعد اشکال می‌تواند دقت پیش‌بینی و سرعت فرآیند یادگیری روش‌های یادگیری ماشین را بهبود بخشد. در این راستا، روش انتخاب ویژگی جدیدی برای بهبود فرآیند پیش‌بینی اشکال نرم‌افزار در این پژوهش ارائه شده است.

در بین روش‌های بسته‌بندی، روش انتخاب روبه‌جلو (FS^2) [۶] عملکرد بسیار خوبی در انتخاب ویژگی‌ها دارد، اما این روش از نظر سرعت بسیار کند است. سرعت پایین این روش به دلیل شروع به کار آن با یک مجموعه خالی از ویژگی‌ها و در ادامه اضافه کردن هر یک از ویژگی‌های مجموعه ویژگی‌ها به‌طور جداگانه به مجموعه اولیه است. ایده پیشنهادی ما، اصلاح مجموعه اولیه ویژگی‌های روش FS با استفاده از تلفیق ویژگی‌های با رتبه بالا در روش‌های فیلتر مختلف است. با توجه به سرعت خوب روش‌های فیلتر، مجموعه پیشنهادی ویژگی‌های اولیه برای شروع در اختیار روش FS قرار می‌گیرد. همچنین مجموعه ویژگی‌های اضافه‌شونده نیز از مابقی ویژگی‌ها تشکیل می‌شود. بنابراین در این پژوهش برای افزایش سرعت و دقت همگرایی روش‌های پایه انتخاب ویژگی، روشی دومرحله‌ای از تلفیق روش‌های فیلتر و بسته‌بندی ارائه شده است در اغلب پژوهش‌های صورت‌پذیرفته در حوزه پیش‌بینی اشکال نرم‌افزار، تنها روش‌های فیلتر مختلف بر دادگان اعمال و نتایج آن‌ها اعلام شده است. در این پژوهش روشی تلفیقی از چند روش فیلتر ارائه شده است که تاکنون در دیگر پژوهش‌ها چنین ترکیبی در تلفیق روش‌های فیلتر ارائه نشده است.

بدون ترکیب روش‌های فیلتر، نتایج پژوهش‌ها نشان‌دهنده آن است که روش‌های فیلتر مختلف بر دادگان مختلف متفاوت عمل کرده‌اند و هر روش فیلتر در برخی دادگان دقت پیش‌بینی بسیار بالا و

می‌کند [۲۱]. به‌طور کلی روش بسته‌بندی کارایی بهتری نسبت به روش فیلتر دارد، اما نسبت به روش فیلتر از نظر محاسباتی پیچیده‌تر است.

در این پژوهش از پنج روش فیلتر شامل روش‌های امتیاز فیشر، شاخص جینی، کروسکال والیس، حداقل افزونگی-حداکثر همبستگی و آزمون T (این روش‌ها در [۱۸] شرح داده شده‌اند) و یک روش بسته‌بندی به نام FS برای ارائه روش انتخاب ویژگی ترکیبی استفاده شده است. علت استفاده از این روش‌ها آن است که این روش‌ها بهترین عملکرد را در پژوهش‌های اخیر در حوزه انتخاب ویژگی در پیش‌بینی اشکال نرم‌افزار [۱۱، ۱۵] داشته‌اند.

۳- کاربرد انتخاب ویژگی در پیش‌بینی اشکال نرم‌افزار

در طول سال‌های اخیر انتخاب ویژگی در کاربردهای بسیاری در حوزه مهندسی نرم‌افزار استفاده شده است [۱۰]. استفاده از این روش‌ها در حوزه پیش‌بینی کیفیت نرم‌افزار و خصوصاً پیش‌بینی اشکال نرم‌افزار نیز روبه‌رشد بوده است [۱۱]. چراکه در این حوزه با این چالش مواجه هستیم که بسیاری از ویژگی‌های استخراج‌شده از ماژول‌های نرم‌افزاری تأثیر منفی در دقت روش‌های یادگیری ماشین دارند و باید کنار گذاشته شوند. پژوهش‌های اخیر که بر تأثیر روش‌های انتخاب ویژگی در پیش‌بینی اشکال نرم‌افزار صورت پذیرفته است را می‌توان به سه دسته تقسیم کرد.

دسته اول پژوهش‌هایی هستند که ایده پیشنهادی آن‌ها تنها بر بررسی تأثیر روش‌های انتخاب ویژگی مختلف بر پیش‌بینی اشکال نرم‌افزار و مقایسه آن‌ها محدود شده است. برای نمونه در [۹] چهار روش مبتنی بر فیلتر و دو روش مبتنی بر بسته‌بندی به‌طور جداگانه و بدون هیچ ایده ترکیبی خاصی برای انتخاب ویژگی‌ها استفاده شده است. در پژوهش مذکور نشان داده شده است که می‌توان با استفاده از کمتر از ده درصد ویژگی‌های اولیه نیز به دقت بالا دست یافت. در [۲۲] نیز، تأثیر ترکیب روش‌های انتخاب ویژگی و یادگیری گروهی^۴ بر روی کارایی این روش‌ها در پیش‌بینی ماژول‌های نرم‌افزاری مستعد اشکال ارزیابی شده است. در این پژوهش ایده پیشنهادی، استفاده از روش یادگیری گروهی بوده است و برای بررسی تأثیر روش‌های انتخاب ویژگی بر کارایی روش یادگیری گروهی نیز از سه روش انتخاب روبه‌جلوی حرصانه (GFS^۵)، روش همبستگی پیرسون و روش امتیاز فیشر استفاده شده است.

دسته دوم پژوهش‌هایی هستند که به‌طور عمده روی جستجوی ویژگی‌های مرتبط با استفاده از ترکیب روش‌های انتخاب ویژگی تمرکز کرده‌اند. برای نمونه در [۱۵] ترکیبی از روش‌های انتخاب ویژگی برای پیش‌بینی اشکال ماژول‌های نرم‌افزاری سیستم مخابراتی بزرگی استفاده شده است. در این پژوهش روشی ترکیبی پیشنهاد شده است که در آن از هفت روش فیلتر و سه روش بسته‌بندی استفاده شده است. در روش ترکیبی ارائه‌شده، ابتدا تمامی ویژگی‌ها توسط روش‌های

پیدا کردن زیرمجموعه‌ای با حداقل اندازه ممکن از ویژگی‌ها هستیم که برای کاربرد موردنظر مناسب باشد. در اغلب موارد تحلیل‌های داده‌ای نظیر طبقه‌بندی بر روی فضای کاسته‌شده نسبت به فضای اصلی بهتر عمل می‌کند [۱۹].

روش‌های انتخاب ویژگی، تلاش می‌کنند تا از میان مجموعه‌ای با N ویژگی و 2^N زیرمجموعه کاندید، بهترین زیرمجموعه را پیدا کنند. در تمام این روش‌ها بر اساس کاربرد و نوع تعریف، زیرمجموعه‌ای به‌عنوان جواب انتخاب می‌شود، که بتواند مقدار یک تابع ارزیابی را بهینه کند. با وجود این که هر روش سعی می‌کند تا بهترین ویژگی‌ها را انتخاب کند، اما با توجه به وسعت جواب‌های ممکن، پیدا کردن جواب بهینه مشکل و در N ‌های متوسط و بزرگ بسیار پرهزینه است [۱۹]. دو دسته‌بندی اصلی از روش‌های انتخاب ویژگی شامل روش‌های فیلتر و بسته‌بندی هستند. در ادامه این بخش این دو روش شرح داده می‌شوند.

۲-۱- روش فیلتر

روش انتخاب ویژگی فیلتر از معیاری کاندید برای امتیازدهی به زیرمجموعه ویژگی‌ها استفاده می‌کند. در این روش، از هیچ روش یادگیری ماشینی استفاده نمی‌شود. به‌عبارتی دیگر، در روش‌های فیلتر، مستقل از الگوریتم یادگیری ماشین، زیرمجموعه‌های ویژگی به‌وسیله مفاهیمی از جمله اطلاعات ویژگی، فاصله، سازگاری، شباهت و دیگر معیارهای آماری انتخاب شده و شایستگی آن‌ها مورد ارزیابی قرار می‌گیرد. در این رویکرد با توجه به معیار کاندید، به هر یک از ویژگی‌ها امتیاز داده می‌شود و بر اساس آن ویژگی‌ها مرتب می‌شوند. این روش‌ها همچنین به‌عنوان پیش‌پردازشی برای روش انتخاب ویژگی بسته‌بندی نیز به‌کار می‌رود [۲۰].

در این رویکرد زیرمجموعه‌ای از ویژگی‌ها بدون استفاده از هیچ الگوریتم یادگیری ماشینی، انتخاب می‌شود بنابراین بسیار سریع عمل می‌کند [۲۰]. اما عیب این روش‌ها این است که نمی‌توانند مشکل افزونگی ویژگی را برطرف کند. افزونگی ویژگی بدان معناست که ویژگی غیرمؤثر بوده و باید کنار گذاشته شود ولی این روش تعداد ویژگی‌های افزونه را مشخص نمی‌کند و به همین دلیل باید تعداد ویژگی‌های مناسب تعیین شوند.

۲-۲- روش بسته‌بندی

به مجموعه روش‌هایی که از تابع ارزیابی مبتنی بر نرخ خطای روش یادگیری استفاده می‌کنند، روش‌های بسته‌بندی یا جعبه سیاه می‌گویند. در این روش هر زیرمجموعه جدید از ویژگی‌ها توسط تابع تولیدکننده ایجاد می‌شود که تولید این زیرمجموعه به استراتژی جستجو وابسته است. آنگاه زیرمجموعه تولیدشده به‌طور مکرر توسط روش یادگیری ماشین ارزیابی می‌شود. تعداد خطاهای مجموعه آزمون، یا همان نرخ خطای روش یادگیری، امتیاز زیرمجموعه را مشخص

توازن کلاس به صورت مجزا محدود شده است. اما در پژوهش مذکور این دو مبحث با یکدیگر ترکیب و به صورت روشی جامع ارائه شده‌اند. در [۱۶] چارچوبی بنام خوشه‌بندی و رتبه‌بندی ویژگی‌ها برای پیش‌بینی ماژول‌های نرم‌افزاری مستعد اشکال پیشنهاد شده است. در چارچوب پیشنهادی آن‌ها برای انتخاب ویژگی از دو مبحث خوشه‌بندی ویژگی‌ها و رتبه‌بندی آن‌ها استفاده شده است. ایشان در روش پیشنهادی ابتدا ویژگی‌های اولیه را بر اساس معیار همبستگی به k خوشه تقسیم می‌کند. سپس مرتبط‌ترین ویژگی‌ها از هر خوشه را مبتنی بر میزان معیار وابستگی آن‌ها انتخاب می‌کند. همچنین در پژوهش خود از روش تقارن بدون قطعیت^{۱۰} به عنوان معیار همبستگی استفاده کرده‌اند. از طرفی دیگر از روش‌های انتخاب ویژگی فیلتری همچون بهره‌وری اطلاعات (IG^{11}) و ReliefF به عنوان معیار وابستگی استفاده شده است. در این پژوهش‌ها مؤثرترین ویژگی‌های نرم‌افزاری جهت بهبود روند پیش‌بینی اشکال نرم‌افزار جستجو و معرفی شده‌اند. نتایج نهایی این پژوهش‌ها نشان‌دهنده تأثیر مثبت روش‌های انتخاب ویژگی در بهبود کارایی پیش‌بینی ماژول‌های نرم‌افزاری مستعد اشکال بوده است.

۴- روش انتخاب ویژگی پیشنهادی

در این پژوهش با توجه به دقت بالاتر روش‌های انتخاب ویژگی بسته‌بندی نسبت به فیلتر و در مقابل سرعت بالای انتخاب ویژگی‌ها توسط روش‌های فیلتر نسبت به بسته‌بندی، از ترکیب این دو روش برای انتخاب بهترین ویژگی‌های مؤثر در بهبود پیش‌بینی اشکال در ماژول‌های نرم‌افزاری استفاده شده است.

ایده پیشنهادی ما، استفاده از روش انتخاب ویژگی دومرحله‌ای برای انتخاب مؤثرترین ویژگی‌ها در فرآیند یادگیری است. روش پیشنهادی در مرحله اول با انتخاب ویژگی‌های اولیه توسط ترکیب وزن دار روش‌های فیلتر موجب بهبود سرعت در انتخاب ویژگی‌ها می‌شود. آنگاه در مرحله دوم با انتخاب ویژگی توسط روش‌های بسته‌بندی با توجه به ویژگی‌های انتخاب‌شده در مرحله اول، بهبود دقت به همراه افزایش سرعت حاصل می‌شود. در ادامه این دو مرحله بیان شده است. هر یک از متغیرهای استفاده‌شده در الگوریتم‌های روش پیشنهادی نیز در جدول ۱ معرفی شده‌اند.

مرحله اول: در روش‌های فیلتر مختلف با توجه به معیار در نظر گرفته‌شده توسط آن‌ها، به هر یک از ویژگی‌ها امتیازی داده می‌شود. امتیاز هر ویژگی در واقع وزن آن ویژگی است که مشخص‌کننده جایگاه ویژگی در لیست مرتب‌شده از ویژگی‌ها توسط روش فیلتر است. ایده پیشنهادی ما در مرحله اول ترکیب وزن دار روش‌های فیلتر است که مطابق الگوریتم ۱ به صورت مراحل زیر است:

۱. وزن‌دهی مجدد به ویژگی‌ها با توجه به اولویت هر ویژگی در لیست مرتب‌شده توسط هر یک از روش‌های فیلتر (اولویت بالاتر در لیست = مقدار وزن بیشتر - سطرهای ۱ تا ۱۰ الگوریتم ۱): در این مرحله

فیلتر مرتب شده‌اند و سپس $\lceil \log_2 n \rceil$ (تعداد کل ویژگی‌ها است) ویژگی اول از لیست مرتب‌شده ویژگی‌ها به عنوان مجموعه کل ویژگی‌ها به روش‌های بسته‌بندی داده شده است. نتایج این پژوهش نشان‌دهنده این است که حذف ۸۵ درصد از ویژگی‌ها تأثیر نامطلوبی روی نتایج نداشته، و در پاره‌ای از موارد نیز حتی موجب بهبود نتایج شده است. ایده پیشنهادی این پژوهش مشابه ایده پیشنهادی ما در روش انتخاب ویژگی دومرحله‌ای است، با این تفاوت که در روش پیشنهادی ما k ویژگی اول از لیست مرتب‌شده ویژگی‌ها (که توسط روش فیلتر ترکیبی مرتب می‌شود) به عنوان مجموعه ویژگی‌های اولیه در روش FS به کار می‌رود و مجموعه ویژگی‌های اضافه‌شونده نیز از مابقی ویژگی‌ها تشکیل می‌شود.

در [۲۳، ۲۴] نیز مطالعه تجربی کاملی بر ۱۷ روش یادگیری گروهی مختلف مبتنی بر روش‌های رتبه‌بندی ویژگی‌های نرم‌افزاری روی ۱۶ مجموعه از دادگان پیش‌بینی اشکال نرم‌افزار، صورت پذیرفته است. روش‌های رتبه‌بندی ویژگی‌ها به روش‌های انتخاب ویژگی مبتنی بر فیلتر اطلاق می‌شود. در این دو پژوهش ایده پیشنهادی روش یادگیری گروهی با محاسبه میانگین نتایج روش‌های فیلتر و انتخاب ویژگی‌های مرتبه بالا بوده است. روش پیشنهادی دو پژوهش فوق مشابه ایده پیشنهادی ما در مرحله اول از روش انتخاب ویژگی دومرحله‌ای (فیلتر ترکیبی وزن دار) است، با این تفاوت که در روش پیشنهادی ما به هر یک از روش‌های فیلتر نیز بر اساس میزان موفقیت آن‌ها در رتبه‌بندی ویژگی‌ها، وزنی نسبت داده شده است. به عبارت دیگر مزیت روش پیشنهادی ما نسبت به ایده پیشنهادی دو پژوهش فوق، اعمال تأثیر هر یک از روش‌های فیلتر بر رتبه هر یک از ویژگی‌های نرم‌افزاری با استفاده از وزن اختصاص یافته به آن روش فیلتر است.

دسته سوم پژوهش‌هایی هستند که به ترکیب روش‌های انتخاب ویژگی با دیگر روش‌های پرکاربرد و موردنیاز در پیش‌بینی اشکال نرم‌افزار پرداخته‌اند. این روش‌ها شامل روش‌های مقابله با مشکل عدم توازن کلاس [۲۵] و خوشه‌بندی دادگان [۱۶] هستند. برای نمونه در [۲۶] ترکیبی از روش‌های بهینه‌سازی فرا-ابتکاری^۶ (روش‌های الگوریتم ژنتیک (GA^۷) و بهینه‌سازی ازدحام ذرات (PSO^۸)) و روش کیسه (Bag^۹) برای بهبود دقت پیش‌بینی اشکال نرم‌افزار ارائه شده است. روش‌های بهینه‌سازی فرا-ابتکاری، برای انتخاب ویژگی و روش Bag برای مقابله با مشکل عدم توازن کلاس استفاده شده است. استفاده از روش‌های بهینه‌سازی فرا-ابتکاری موجب افزایش توانایی یافتن راه‌حل‌هایی با کیفیت بالا در بازه زمانی معقول می‌شود. هدف پژوهش مذکور از انتخاب روش‌های فوق برای انتخاب ویژگی بهبود دقت پیش‌بینی بوده است. به عبارت دیگر فاکتور سرعت از اهداف اصلی آن نبوده است، بنابراین از روش‌های انتخاب ویژگی فیلتر در فرآیند انتخاب ویژگی استفاده نشده است. اکثر پژوهش‌های حوزه پیش‌بینی اشکال، ایده پیشنهادی آن‌ها به انتخاب ویژگی و یا حل مشکل عدم

M روش فیلتر داشته باشیم که $fm_j.W$ بیانگر وزن روش فیلتر نام است. وزن نهایی با رابطه (۱) محاسبه می‌شود.

$$x.NewW = \frac{\sum_{j=1}^M fm_j.W \times x.W}{M} \quad (1)$$

مرحله ۳ برای اعمال تأثیر میزان موفقیت روش‌های فیلتر در وزن دهی به ویژگی‌ها است. این روش به‌عنوان مرحله اول از روش انتخاب ویژگی دومرحله‌ای استفاده می‌شود به همین دلیل در سطرهای ۳۷ تا ۴۷ الگوریتم ۱ میزان k بهینه برای این روش محاسبه شده است، چراکه یکی از پارامترهای ورودی مرحله دوم است.

مرحله دوم: این مرحله با توجه به اولویت‌بندی ویژگی‌ها توسط روش فیلتر ترکیبی در مرحله اول آغاز می‌شود. در این مرحله مطابق الگوریتم ۲، از روش انتخاب ویژگی بسته‌بندی FS برای کاهش ویژگی‌ها و بهبود دقت پیش‌بینی اشکال در ماژول‌های نرم‌افزاری استفاده شده است. در روش FS، k ویژگی اول از روش فیلتر ترکیبی به‌عنوان مجموعه اولیه و باقی ویژگی‌ها به‌عنوان مجموعه اضافه شونده در نظر گرفته شده است. بنابراین ما این روش را توسعه داده‌ایم. مراحل انتخاب ویژگی پیشنهادی در مرحله دوم به‌صورت الگوریتم ۲ است. همچنین مراحل انتخاب ویژگی توسط روش FS معمولی نیز به‌صورت الگوریتم ۳ (بازنویسی مجدد الگوریتم ۱ از [۲۷]) است.

بر اساس لیست مرتب‌شده از ویژگی‌ها که روش فیلتر به‌عنوان خروجی برمی‌گرداند، وزندهی مجدد انجام می‌شود و به ویژگی‌های مرتبط‌تر (و دارای اولویت بالاتر در لیست) وزن (w) بیشتری داده می‌شود.

۲. محاسبه وزن برای هر یک از روش‌های فیلتر بر اساس میزان موفقیت آن‌ها در اولویت‌بندی مؤثر ویژگی‌ها با استفاده از دادگان آموزش و اعتبارسنجی (سطرهای ۱۱ تا ۳۰ الگوریتم ۱): در این مرحله با هر یک از روش‌های فیلتر ویژگی‌های دادگان آموزش را اولویت‌بندی می‌کنیم و با بخش اعتبارسنجی میزان دقت هر یک از روش‌های فیلتر را ارزیابی می‌کنیم. آنگاه روش‌ها را بر اساس میزان موفقیتشان مرتب کرده و به روشی که کمترین دقت را دارد وزن 0.1 اختصاص می‌یابد و به دیگر روش‌ها بر اساس میزان اختلاف دقت آن‌ها نسبت به دقت روش قبلی وزنی معادل «میزان اختلاف تقسیم بر ۱۰ و به‌اضافه وزن روش قبلی» اختصاص می‌دهیم. برای نمونه به روشی که میزان دقت آن ۵ درصد بیشتر از روش با وزن 0.1 است، وزن 0.16 تعلق می‌گیرد.

۳. وزن دهی نهایی به ویژگی با ضرب وزن آن در وزن هر یک از روش‌های فیلتر و میانگین‌گیری از وزن جدید ویژگی با محاسبه حاصل جمع وزن‌های ویژگی در هر روش فیلتر و تقسیم آن به تعداد کل روش‌های فیلتر (سطرهای ۳۱ تا ۳۶ الگوریتم ۱): اگر ویژگی را با x نشان دهیم که دارای وزن اولیه W است و همچنین

جدول ۱: متغیرهای استفاده‌شده در الگوریتم‌های روش پیشنهادی

نماد متغیر	شرح متغیر	نماد متغیر	شرح متغیر
X^{train}	مجموعه دادگان آموزشی	accuracy	میزان دقت
x^i	یک نمونه آموزشی	best	بهترین میزان دقت در هر تکرار الگوریتم انتخاب ویژگی
r^i	برچسب نمونه آموزشی (مستعد اشکال / غیر مستعد اشکال)	$fm_i.bestAcc$	میزان دقت روش فیلتر نام
N_t	تعداد کل نمونه آموزشی	FilterList	لیست مرتب‌شده روش‌های فیلتر بر اساس میزان دقت آن‌ها
D	تعداد ویژگی‌های نمونه آموزشی	$fm_i.W$	وزن روش فیلتر نام
$X^{validate}$	مجموعه دادگان آزمون	d	میزان اختلاف دقت روش فیلتر نام با روش فیلتر 1-نام
x^v	یک نمونه آزمون	$x_i^v.NewW$	میزان وزن جدید ویژگی نام نمونه آموزشی با روش فیلتر ترکیبی وزن‌دار پیشنهادی
r^v	برچسب نمونه آزمون (مستعد اشکال / غیر مستعد اشکال)	newFList	لیست ویژگی‌های مرتب‌شده توسط روش فیلتر ترکیبی وزن‌دار پیشنهادی
FM	مجموعه روش‌های فیلتر	optimalk	تعداد ویژگی‌های بهینه برای روش فیلتر ترکیبی وزن‌دار پیشنهادی
N_v	تعداد کل نمونه آزمون	S	لیست کل ویژگی‌های دادگان
fm_i	روش فیلتر نام	IS	لیست ویژگی‌های اولیه در روش FS
M	تعداد کل روش‌های فیلتر	AS	لیست ویژگی‌های اضافه‌شونده در روش FS
$FList_i$	لیست ویژگی‌های مرتب‌شده توسط روش فیلتر نام	FSFList	لیست ویژگی‌های نهایی انتخاب‌شده توسط روش FS
w	متغیری برای نگهداری موقت تعداد ویژگی‌ها	EM	میزان معیار ارزیابی (دقت)
$fnum$	شماره ویژگی	Max	متغیری برای نگهداری تعداد کل ویژگی‌های اضافه‌شونده در FS
$x_{fnum}^i.W$	وزن ویژگی با شماره $fnum$ از نمونه آموزشی	s	ویژگی انتخاب‌شده از مجموعه ویژگی‌های اضافه‌شونده
train	متغیری برای نگهداری موقت نمونه‌های آموزشی با ویژگی‌های انتخاب‌شده توسط روش فیلتر	best-feature	متغیری برای نگهداری موقت بهترین ویژگی منتخب FS در هر تکرار الگوریتم
validate	متغیری برای نگهداری موقت نمونه‌های آزمون با ویژگی‌های انتخاب‌شده توسط روش فیلتر	n	متغیری برای نگهداری موقت تعداد کل ویژگی‌های باقی‌مانده از مجموعه ویژگی‌های اضافه‌شونده در روش FS

الگوریتم ۱: روش فیلتر ترکیبی وزن‌دار پیشنهادی (مرحله اول از روش انتخاب ویژگی دومرحله‌ای)

Fused weighted filter method (X^{train}, X^{validate}, FM)

Input:

$X^{train} = \{x^t, r^t\}_{t=1}^{Nt}$ // x^t is a train sample, $r^t \in \{nfp, fp\}$, Nt is the number of train samples, $x^t \in [x^t_1, x^t_2, x^t_3, \dots, x^t_D]$ and D is the number of features.

$X^{validate} = \{x^v, r^v\}_{v=1}^{Nv}$ // x^v is a validate sample, $r^v \in \{nfp, fp\}$, Nv is the number of validate samples, $x^v \in [x^v_1, x^v_2, x^v_3, \dots, x^v_D]$ and D is the number of features.

$FM = \{fm_i\}_{i=1}^M$ // fm is Filter method and M is the number of Filter methods.

Method:

```

1:   for i=1: M
2:       FListi = fmi (xt); //FListi is ordered List of features by fmi.
3:       w=D;
4:       // Assignment new weight to features for each of filter methods.
5:       for j=1: D
6:           fnum=FListi(j); //fnum is feature's number.
7:           xtfnum.W = w;
8:           w=w-1;
9:       end
10:  end
11:  // calculate weight for each of filter methods.
12:  for i=1: M
13:      FList = fmi (xt); // FList is ordered List of features by fmi.
14:      for j=1: size(FList)
15:          train = train + FList.xtj;
16:          validate = validate + FList.xvj;
17:          accuracy = classifier (train, validate);
18:          if accuracy>best result in this iteration
19:              best=accuracy;
20:          end
21:      end
22:      fmi.bestAcc= best;
23:  end
24:  FilterList= sort (fm.bestAcc, 'ascend'); // FilterList is ordered list of Filter methods.
25:  i=FilterList(1);
26:  fmi.W=0.1;
27:  for i=2: M
28:      d= fmi.bestAcc - fmi-1.bestAcc;
29:      fmi.W= fmi-1.W + (d/10);
30:  end
31:  // Assignment new weight to features by WF method.
32:  for i=1: D
33:      xti.NewW =  $\frac{\sum_{j=1}^M fm_j.W \times x^t_j.W}{M}$ ;
34:      newFList.Add ( xti );
35:  end
36:  [newFList] = sort (newFList.NewW, 'descend');
37:  // calculate optimal k for feature selection by WF method.
38:  for j=1: size(newFList)
39:      train = train + newFList.xtj;
40:      validate = validate + newFList.xvj;
41:      accuracy = classifier (train, validate);
42:      if accuracy>best result in this iteration
43:          best=accuracy;
44:          k=j;
45:      end
46:  end
47:  optimalk= k;

```

Output:

newFList // the ordered subset of the features by WF method.
 optimalk // the optimal k for feature selection by WF method.

الگوریتم ۲: روش انتخاب روبه‌جلو توسعه‌یافته پیشنهادی (مرحله دوم از روش انتخاب ویژگی دومرحله‌ای)

در دو پاراگراف آخر از صفحه ششم مقاله ۲ بار به این الگوریتم اشاره شده است. **Enhanced forward feature selection method (X, FM)**

Input:

$X = \{x^t, r^t\}_{t=1}^N$ // x^t is a train sample, $r^t \in \{\text{nfp}, \text{fp}\}$, N is the number of samples, $x^t \in [x_1^t, x_2^t, x_3^t, \dots, x_D^t]$, and D is the number of features.

$FM = \{fm_i\}_{i=1}^M$ // fm_i is a known Filter method and M is the number of Filter methods.

Method:

2: $[X^{\text{train}} X^{\text{validate}}] = 10\text{-fold cross validation}(X)$;

3: $[FList\ k] = \text{Fused weighted filter method}(X^{\text{train}}, X^{\text{validate}}, FM)$;

4: // Enhanced forward feature selection method:

4: $S = FList$; // $FList$ is $[x_1^t, x_2^t, x_3^t, \dots, x_D^t]$.

5: $IS = [x_1^t, x_2^t, x_3^t, \dots, x_k^t]$; // IS is initial subset for enhanced forward feature selection method.

6: $AS = [x_{k+1}^t, x_{k+2}^t, x_{k+3}^t, \dots, x_D^t]$; // AS is $S-IS$ and it is additive subset for enhanced forward feature selection method.

7: $[FSFList, EM] = \text{Forward selection method}(X^{\text{train}}, X^{\text{validate}}, IS, AS)$

Output:

$FSFList$ // the optimal subset of the features by forward selection method.

EM // the value of evaluation measures by forward selection method.

الگوریتم ۳: روش FS

Forward selection method ($X^{\text{train}}, X^{\text{validate}}, IS, AS$)**Input:**

$X^{\text{train}} = \{x^t, r^t\}_{t=1}^{Nt}$ // x^t is a train sample, $r^t \in \{\text{nfp}, \text{fp}\}$, Nt is the number of train samples, $x^t \in [x_1^t, x_2^t, x_3^t, \dots, x_D^t]$, and D is the number of features.

$X^{\text{validate}} = \{x^v, r^v\}_{v=1}^{Nv}$ // x^v is a validate sample, $r^v \in \{\text{nfp}, \text{fp}\}$, Nv is the number of validate samples, $x^v \in [x_1^v, x_2^v, x_3^v, \dots, x_D^v]$, and D is the number of features.

IS // initial subset for forward feature selection method.

AS // additive subset for forward feature selection method.

Method:

1: $n = 1$;

2: $Max = \text{Size}(AS)$;

3: while ($n \leq Max$)

4: while ($s = \text{selected next element of } AS$)

5: $S = IS \cup s$;

6: $[accuracy] = \text{Classifier}(X^{\text{train}}, X^{\text{validate}}, S)$;

7: If $accuracy > \text{Best result in this iteration}$

8: $Best = accuracy$;

9: $best\text{-feature} = s$;

10: end

11: end

12: if $Best > \text{previous Best result}$

13: $IS = IS \cup \text{best-feature}$;

14: $AS = AS - \text{best-feature}$;

15: $EM = Best$;

16: else

17: break;

18: end

19: $n = n + 1$;

20: end

Output:

IS // the optimal subset of the features

EM // the best value of evaluation measures by features in IS

۵- دادگان‌های ارزیابی مدل پیشنهادی

مقادیر غیر صفر نشان‌دهنده مستعد اشکال بودن آن است. جدول ۲ اطلاعات این دادگان را نشان می‌دهد. تعداد نمونه‌ها نشان‌دهنده تعداد ماژول‌های نرم‌افزاری است. هر نمونه بیانگر ماژولی است که ویژگی‌های مختلفی از آن استخراج شده است. ویژگی‌های هر یک از دادگان نیز در جدول مذکور با علامت × نشان داده شده است [۲۸].

جهت بررسی مؤثر بودن روش‌های پیشنهادی در پیش‌بینی اشکال نرم‌افزار، از ۸ دادگان مخصوص پیش‌بینی اشکال نرم‌افزار استفاده شده است. در این دادگان‌ها، آخرین ویژگی برچسب کلاس هر نمونه است. در این ویژگی، مقدار صفر نشان‌دهنده غیرمستعد اشکال بودن ماژول و

جدول ۲: اطلاعات دادگان پیش‌بینی اشکال نرم‌افزار

PC3	PC2	PC1	MW1	MC2	MC1	KC3	CM1	معادل فارسی	نام کامل ویژگی	نماد ویژگی
×	×	×	×	×	×	×	×	تعداد کل خطوط کد	Line count of code	X ₁
×	×	×	×	×	×	×	×	پیچیدگی نمودار جریان	McCable's cyclomatic complexity	X ₂
×	×	×	×	×	×	×	×	پیچیدگی تجزیه	McCable's essential complexity	X ₃
×	×	×	×	×	×	×	×	پیچیدگی طراحی	McCable's design complexity	X ₄
×	×	×	×	×	×	×	×	تعداد کل عملگرها	Total number of operators	X ₅
×	×	×	×	×	×	×	×	تعداد کل عملوندها	Total number of operands	X ₆
×	×	×	×	×	×	×	×	تعداد کل عملگرهای یکتا	Number of unique operators	X ₇
×	×	×	×	×	×	×	×	تعداد کل عملوندهای یکتا	Number of unique operands	X ₈
×	×	×	×	×	×	×	×	تعداد کل عملگرها و عملوندهای یکتا	Number of unique operators and operands	X ₉
×	×	×	×	×	×	×	×	حجم	Halstead's volume	X ₁₀
×	×	×	×	×	×	×	×	دشواری	Halstead's difficult	X ₁₁
×	×	×	×	×	×	×	×	طول ماژول	Halstead's length	X ₁₂
×	×	×	×	×	×	×	×	محتوا	Halstead's content	X ₁₃
×	×	×	×	×	×	×	×	میزان تلاش	Halstead's effort	X ₁₄
×	×	×	×	×	×	×	×	میزان برآورد خطا	Halstead's error estimate	X ₁₅
×	×	×	×	×	×	×	×	میزان زمان برنامه‌نویسی	Halstead's programing time	X ₁₆
×	×	×	×	×	×	×	×	تعداد خطوط خالی ماژول	Number of blank lines	X ₁₇
×	×	×	×	×	×	×	×	تعداد خطوط توضیحات از ماژول	Number of comment-only lines	X ₁₈
		×					×	تعداد خطوطی که تنها شامل کد هستند	Number of code-only lines	X ₁₉
×	×	×	×	×	×	×	×	تعداد خطوط توضیحات و کد	Number of lines with both code and comments	X ₂₀
×	×	×	×	×	×	×	×	تعداد کل نمودارهای جریان	Number of branches	X ₂₁
×	×		×	×	×	×		تعداد دستورهای شرطی	Number of condition	X ₂₂
×	×		×	×	×	×		فراخوانی دوجزئی	Call pairs	X ₂₃
×	×		×	×	×	×		تراکم مسیرهای نمودار جریان	Cyclomatic density	X ₂₄
×	×		×	×	×	×		تعداد دستورهای تصمیم‌گیری	Number of decision	X ₂₅
×	×		×	×	×	×		تراکم مسیرهای دستورهای تصمیم‌گیری	Decision density	X ₂₆
×	×		×	×	×	×		تراکم طراحی	Design density	X ₂₇
×	×		×	×	×	×		تعداد لبه‌ها یا مرزها	Number of edge	X ₂₈
×	×		×	×	×	×		تراکم ضروری	Essential density	X ₂₉
×	×		×	×	×	×		تعداد خط‌های قابل‌اجرا	LOC executable	X ₃₀
×	×		×	×	×	×		تعداد پارامترهای ورودی	Number of parameter	X ₃₁
				×	×	×		پیچیدگی داده‌های سراسری	Global data complexity	X ₃₂
				×	×	×		تراکم داده‌های سراسری	Global data density	X ₃₃
×	×		×	×	×	×		سطح	Halstead's level	X ₃₄
×	×		×	×	×	×		شدت تعمیر و نگهداری	Maintenance severity	X ₃₅
×	×		×	×	×	×		تعداد شرط‌های اصلاح‌شده	Number of modified condition	X ₃₆
×	×		×	×	×	×		تعداد شرط‌های چندگانه	Number of multiple condition	X ₃₇
×	×		×	×	×	×		تعداد گره‌ها	Number of node	X ₃₈
×	×		×	×	×	×		میزان نرمال‌شده تعداد مسیرهای مستقل خطی در نمودار جریان ماژول	Normalized cyclomatic complexity	X ₃₉
×	×		×	×	×	×		میزان درصد توضیحات در کد	Percent comments	X ₄₀
۳۷	۳۷	۲۱	۳۷	۳۹	۳۹	۳۹	۲۱		تعداد ویژگی‌ها	
۷۵۳	۷۷۸	۱۱۰۹	۴۰۳	۱۶۱	۸۹۶	۴۵۸	۴۹۸		تعداد نمونه‌ها	

۶- معیارهای ارزیابی مدل پیشنهادی

از جمله عوامل مؤثر در مقایسه مدل با دیگر مدل‌ها، سنجش آن با معیار ارزیابی مناسب است. در این پژوهش جهت ارزیابی دقت طبقه‌بندها در تشخیص ماژول‌های نرم‌افزاری مستعد اشکال با توجه به معیارهای استفاده‌شده در پژوهش‌هایی که با نتایج روش‌های پیشنهادی ما مورد مقایسه قرار گرفته‌اند، از معیار سطح زیر منحنی مشخصه عملکرد سیستم (AUC^{12}) استفاده شده است. معیار AUC به‌طور گسترده برای اندازه‌گیری میزان کارایی مدل طبقه‌بندی استفاده می‌شود. این معیار را می‌توان به‌عنوان یک توصیف آماری استفاده کرد، تا تخمین زده شود که احتمال این که یک مدل پیش‌بینی، یک نمونه دارای اشکال را بیشتر از یک نمونه فاقد اشکال، تشخیص دهد، چه مقدار است [۲۹].

در این معیار، از چهار معیار شامل معیارهای مثبت واقعی (TP^{13})، مثبت کاذب (FP^{14})، منفی واقعی (TN^{15})، منفی کاذب (FN^{16})، نرخ مثبت واقعی (TP_r^{17}) و نرخ مثبت کاذب (FP_r^{18}) استفاده شده است. TP بیانگر تعداد ماژول‌های دارای اشکال که به‌درستی مستعد اشکال پیش‌بینی شده‌اند، است. FP بیانگر تعداد ماژول‌های فاقد اشکال است که به‌اشتباه مستعد اشکال پیش‌بینی شده‌اند. TN بیانگر تعداد ماژول‌های فاقد اشکال که به‌درستی غیرمستعد اشکال پیش‌بینی شده‌اند، است. FN نیز تعداد ماژول‌های دارای اشکال که به‌اشتباه غیرمستعد اشکال پیش‌بینی شده‌اند، است. معیار TP_r بیانگر نرخ تشخیص صحیح دسته مثبت توسط طبقه‌بند است. معیار FP_r نرخ تشخیص اشتباه دسته منفی را بیان می‌کند [۲۹]. دو معیار TP_r و FP_r به‌ترتیب با رابطه‌های (۲) و (۳) محاسبه می‌شوند. منظور از دسته مثبت، نمونه‌های دارای برجسب کلاس غیرصفر هستند که همان ماژول‌های مستعد اشکال می‌باشند. همچنین منظور از دسته منفی نمونه‌های دارای برجسب کلاس صفر یعنی ماژول‌های غیر مستعد اشکال هستند. مقدار معیار AUC با رابطه (۴) محاسبه می‌شود [۳۰].

$$TP_r = \frac{TP}{TP + FN} \quad (2)$$

$$FP_r = \frac{FP}{FP + TN} \quad (3)$$

$$AUC = \frac{1 + TP_r - FP_r}{2} \quad (4)$$

۷- نتایج تجربی آزمایش‌ها

تمامی دادگان در این پژوهش ابتدا به دو بخش دادگان آموزش و دادگان آزمون تقسیم شده‌اند. تقسیم‌بندی به‌صورت، ۹۰ درصد برای آموزش و ۱۰ درصد برای آزمون (ارزیابی A1) و ۱۰ درصد برای آموزش و ۹۰ درصد برای آزمون (ارزیابی A2) انجام شده است. آنگاه بخش آموزش هم، خود به دو بخش آموزش و اعتبارسنجی تقسیم شده است که برای تعیین پارامترهای مدل پیشنهادی مورد استفاده قرار

گرفته است. آنچه مسلم است این است که نتیجه حاصل از ارزیابی دادگانی که در آموزش استفاده شده‌اند به‌عنوان نتیجه ارزیابی کلی مدنظر قرار نمی‌گیرد، بلکه آنچه به‌عنوان نتیجه ارزیابی کارایی روش مدنظر است، دقت آن روش روی پیش‌بینی ماژول‌های نرم‌افزاری مستعد اشکالی که در بخش دادگان آزمون قرار دارند، است. برای این تقسیم‌بندی بخش آموزش به دو بخش آموزش و اعتبارسنجی از روش اعتبارسنجی متقابل k بخشی 19 با $k=10$ استفاده شده است. همچنین برای طبقه‌بندی ماژول‌های نرم‌افزاری به دو دسته مستعد اشکال و غیرمستعد اشکال در این پژوهش از روش ماشین بردار پشتیبان دوقلو هموارشده با حاشیه پارامتریک ($STPM SVM^{20}$) [۱۷] استفاده شده است. علت استفاده از روش مذکور به‌عنوان طبقه‌بند، برتری این روش نسبت به دیگر طبقه‌بندها در پیش‌بینی اشکال نرم‌افزار در پژوهش [۳۱] است.

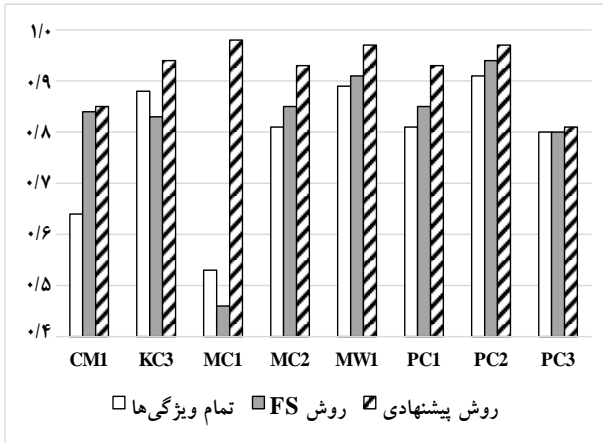
مرحله دوم از روش انتخاب ویژگی دومرحله‌ای، با ویژگی‌های به‌دست‌آمده توسط روش فیلتر ترکیبی وزن‌دار در مرحله اول، آغاز می‌شود. در این مرحله از روش‌های انتخاب ویژگی بسته‌بندی برای کاهش ویژگی‌ها و بهبود دقت پیش‌بینی ماژول‌های نرم‌افزاری مستعد اشکال استفاده شده است. روش بسته‌بندی استفاده‌شده در این مرحله روش FS است که پارامترهای ورودی آن مطابق الگوریتم ۲ اصلاح شده است و اصطلاحاً توسعه داده شده است.

برای روش فیلتر ترکیبی پیشنهادشده، میانگین تعداد ویژگی‌های انتخاب‌شده از هر یک از دادگان حدود ۵۶ درصد است. بنابراین حدود ۵۶ درصد از ویژگی‌ها (که برترین ویژگی‌ها از نظر روش فیلتر ترکیبی وزن‌دار پیشنهادی هستند) ثابت در نظر گرفته می‌شود و ۴۴ درصد از ویژگی‌های باقی‌مانده توسط روش FS موردبررسی قرار می‌گیرد که این مسئله علاوه بر بهبود دقت، موجب نصف شدن زمان اجرای روش FS می‌شود. بنابراین می‌توان نتیجه گرفت که در پروژه‌های نرم‌افزاری واقعی که تعداد ویژگی‌ها و ماژول‌های نرم‌افزاری بسیار بیشتر از دادگان موردآزمایش قرار گرفته در این پژوهش هستند، الگوریتم پیشنهادی می‌تواند با انتخاب ویژگی‌های برتر، علاوه بر بهبود دقت پیش‌بینی، موجب کاهش زمان و حدوداً دو برابر شدن سرعت نیز شود.

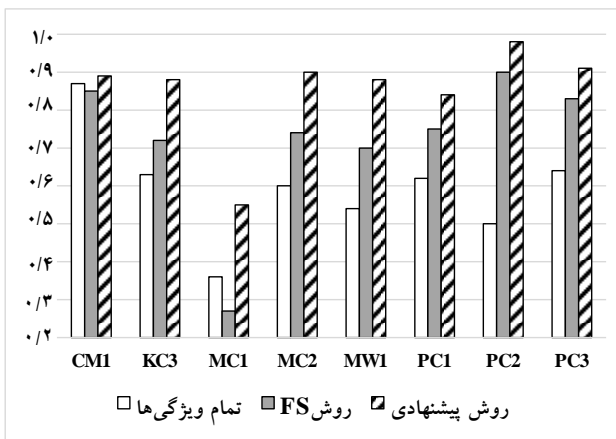
در آزمایش روش پیشنهادی علاوه بر ارزیابی دقت، اقدام به اندازه‌گیری زمان اجرای الگوریتم پیشنهادی (شامل روش فیلتر ترکیبی وزن‌دار و روش FS توسعه‌داده‌شده) و روش FS معمولی نموده‌ایم. اندازه‌گیری زمان با استفاده از توابع Tic و Toc نرم‌افزار متلب انجام شده است. تمامی الگوریتم‌های استفاده‌شده در روش پیشنهادی در محیط متلب و روی یک کامپیوتر شخصی با پردازنده چهار هسته‌ای، RAM معادل ۴ گیگابایت و با ویندوز ۷ پیاده‌سازی شده است.

نتایج نهایی روش انتخاب ویژگی ترکیبی دومرحله‌ای و همچنین نتایج روش FS معمولی و نتایج عدم استفاده از این روش (استفاده از تمام ویژگی‌ها) در جدول‌های ۳ و ۴ به‌صورت دقت و زمان اجرا (برحسب ثانیه) نشان داده شده است. میزان بهبود دقت روش

سریع‌تر عمل کرده است. همچنین مطابق جدول ۴ مشاهده می‌شود که روش پیشنهادی در ارزیابی A2 نسبت به روش FS معمولی به‌طور میانگین ۰/۴۷ ثانیه سریع‌تر عمل کرده است. منحنی میزان AUC روش پیشنهادی بر ۸ دادگان ارزیابی ناسا برای ارزیابی‌های A1 و A2 به ترتیب در نمودارهای شکل ۱ و شکل ۲ نشان داده شده است.



شکل ۱: نمودار منحنی میزان AUC روش پیشنهادی در مقایسه با روش FS و استفاده از تمام ویژگی‌ها با ارزیابی A1



شکل ۲: نمودار منحنی میزان AUC روش پیشنهادی در مقایسه با روش FS و استفاده از تمام ویژگی‌ها با ارزیابی A2

با توجه به نمودارهای شکل ۱ و شکل ۲ مشاهده می‌شود که روش پیشنهادی میزان AUC قابل قبولی برای پیش‌بینی اشکال در این دادگان داشته است و در تمامی دادگان میزان دقت با این روش بهبود یافته است.

بررسی ویژگی‌های انتخاب‌شده از دادگان توسط روش پیشنهادی بیانگر آن است که پنج ویژگی بهترین و مؤثرترین ویژگی‌های نرم‌افزاری برای بهبود پیش‌بینی اشکال نرم‌افزار با روش انتخاب ویژگی پیشنهادی به‌طور مشترک در هر دو ارزیابی A1 و A2 هستند. این ویژگی‌ها در جدول ۵ آورده شده‌اند. همچنین با بررسی دقیق‌تر نتایج مشاهده می‌شود که سه ویژگی نیز غیرمؤثرترین ویژگی‌های نرم‌افزاری در پیش‌بینی اشکال نرم‌افزار با روش انتخاب ویژگی پیشنهادی به‌طور

پیشنهادی نسبت به استفاده از تمام ویژگی‌ها و میزان افزایش سرعت روش پیشنهادی نسبت به روش FS معمولی نیز در جدول‌های مذکور نشان داده شده است.

جدول ۳: نتایج روش انتخاب ویژگی پیشنهادی در مقایسه با نتایج روش FS معمولی و استفاده از تمام ویژگی‌ها با ارزیابی A1

میزان افزایش سرعت	میزان بهبود دقت	روش پیشنهادی		روش FS معمولی		تمام ویژگی‌ها	دادگان
		زمان	دقت	زمان	دقت		
+۰/۵۹	+۰/۲۱	۰/۵۲	۰/۸۵	۱/۱۱	۰/۸۴	۰/۶۴	CM1
+۰/۵۴	+۰/۰۶	۰/۲۴	۰/۹۴	۰/۷۸	۰/۸۳	۰/۸۸	KC3
+۰/۵	+۰/۴۵	۱/۴۷	۰/۹۸	۱/۹۷	۰/۴۶	۰/۵۳	MC1
+۰/۶۷	+۰/۱۲	۰/۱۹	۰/۹۳	۰/۸۶	۰/۸۵	۰/۸۱	MC2
+۰/۴۹	+۰/۰۸	۰/۲۰	۰/۹۷	۰/۶۹	۰/۹۱	۰/۸۹	MW1
+۰/۵۱	+۰/۱۲	۰/۸۶	۰/۹۳	۱/۳۷	۰/۸۵	۰/۸۱	PC1
+۱/۰۸	+۰/۰۶	۰/۴۶	۰/۹۷	۱/۵۴	۰/۹۴	۰/۹۱	PC2
+۱/۶۹	+۰/۰۱	۰/۳۷	۰/۸۱	۲/۰۶	۰/۸۰	۰/۸۰	PC3
+۰/۷۶	+۰/۱۴	۰/۵۴	۰/۹۲	۱/۳۰	۰/۸۱	۰/۷۸	میانگین

جدول ۴: نتایج روش انتخاب ویژگی پیشنهادی در مقایسه با نتایج روش FS معمولی و استفاده از تمام ویژگی‌ها با ارزیابی A2

میزان افزایش سرعت	میزان بهبود دقت	روش پیشنهادی		روش FS معمولی		تمام ویژگی‌ها	دادگان
		زمان	دقت	زمان	دقت		
+۰/۹	+۰/۰۴	۰/۵۱	۰/۸۹	۱/۴۱	۰/۸۵	۰/۸۷	CM1
+۰/۷۲	+۰/۱۶	۰/۰۶	۰/۸۸	۰/۷۸	۰/۷۲	۰/۶۳	KC3
+۰/۳۳	+۰/۲۸	۰/۶۱	۰/۵۵	۰/۹۴	۰/۲۷	۰/۳۶	MC1
+۰/۶۵	+۰/۱۶	۰/۱۰	۰/۹۰	۰/۷۵	۰/۷۴	۰/۶۰	MC2
+۰/۳۸	+۰/۱۸	۰/۲۲	۰/۸۸	۰/۶۰	۰/۷۰	۰/۵۴	MW1
+۰/۰۳	+۰/۱۲	۰/۸۴	۰/۸۴	۰/۸۷	۰/۷۵	۰/۶۲	PC1
+۰/۱۱	+۰/۰۸	۰/۰۹	۰/۹۸	۰/۲۰	۰/۹۰	۰/۵۰	PC2
+۰/۶۱	+۰/۰۸	۰/۱۴	۰/۹۱	۰/۷۵	۰/۸۳	۰/۶۴	PC3
+۰/۴۷	+۰/۱۴	۰/۳۲	۰/۸۵	۰/۷۰	۰/۷۲	۰/۵۹	میانگین

مطابق جدول ۳ مشاهده می‌شود که در تمامی دادگان میزان دقت بهبود یافته است. برای نمونه روش پیشنهادی در دادگان CMI با میزان دقت، ۰/۸۵، مازول‌های نرم‌افزاری مستعد اشکال را به‌درستی پیش‌بینی کرده است که بیشتر از میزان دقت به‌دست‌آمده با تمام ویژگی‌ها (۰/۶۴) است. نکته قابل توجه آن است که روش پیشنهادی به‌طور میانگین ۱۴ درصد بهبود دقت داشته است. همچنین این نکته نیز قابل توجه است که در روش پیشنهادی حدود ۵۶ درصد از ویژگی‌ها در مرحله اول توسط روش فیلتر انتخاب شده‌اند و مابقی ویژگی‌ها توسط روش FS موردبررسی قرار گرفته‌اند، بنابراین می‌توان نتیجه گرفت که با کاهش ۵۰ درصدی پیچیدگی زمانی نسبت به، استفاده از روش FS معمولی، روش پیشنهادی سرعت اجرای حدوداً دو برابر داشته است. مطابق جدول ۳ مشاهده می‌شود که روش پیشنهادی نسبت به روش FS معمولی به‌طور میانگین ۰/۷۶ ثانیه در ارزیابی A1

ویژگی برای انتخاب مؤثرترین ویژگی‌های ماژول‌های نرم‌افزاری ارائه شدند. روش بسته‌بندی به‌علت استفاده از یک الگوریتم یادگیری ماشین دارای سرعت پایینی است اما این روش دقت بالایی دارد. در مقابل روش فیلتر سرعت عمل بالایی دارد ولی دقت چندان مناسبی را ارائه نمی‌دهد. با ترکیب نقاط قوت این دو روش یعنی عملکرد خوب روش انتخاب ویژگی بسته‌بندی و همچنین سرعت بالای روش انتخاب ویژگی فیلتر، ما به روش ترکیبی مناسبی برای انتخاب ویژگی‌ها دست یافته‌ایم.

در این پژوهش ما روش FS که یکی از بهترین روش‌های انتخاب ویژگی بسته‌بندی است را توسعه داده‌ایم. روش FS برای شروع به کار خود از یک مجموعه خالی به‌عنوان مجموعه ویژگی‌های اولیه استفاده می‌کند. ما در این پژوهش این مجموعه اولیه را با ویژگی‌های به‌دست‌آمده از ترکیب روش‌های فیلتر اصلاح کرده‌ایم. سپس این روش را روی دادگان پیش‌بینی اشکال نرم‌افزار ناسا مورد ارزیابی قرار دادیم. نتایج حاکی از مؤثر بودن روش مورد استفاده می‌باشند. بنابراین روش پیشنهادی ما توانسته است، بهترین ویژگی‌ها برای بهبود دقت پیش‌بینی اشکال را با سرعت بیشتری بیابد. در این پژوهش ما از پنج روش فیلتر و یک روش بسته‌بندی استفاده کردیم. ما قصد داریم در آینده روی ترکیبات دیگری از روش‌های فیلتر و بسته‌بندی کار کنیم.

مراجع

- [1] I. Gondra, "Applying machine learning to software fault-proneness prediction," *J. Syst. Softw.*, vol. 81, no. 2, pp. 186-195, 2008.
- [2] علیقارداشی و زارع چاهوکی، «افزایش دقت در پیش‌بینی اشکال نرم‌افزار با استفاده از ماشین بردار پشتیبان دوقلو با طرح حداقل مربعات بازگشتی»، هشتمین کنفرانس داده‌کاوی ایران، تهران، ۱۳۹۳.
- [3] E. Rashid, S. Patnaik, and A. Usmani, "Machine learning and its application in software fault prediction with similarity measures," *Comput. Vis. Robot.*, Springer India, vol. 332, pp. 37-45, 2015.
- [4] C. Catal, "Software fault prediction: A literature review and current trends," *Expert Syst. Appl.*, vol. 38, no. 4, pp. 4626-4636, 2011.
- [5] T. Menzies, J. Greenwald, and A. Frank, "Data mining static code attributes to learn defect predictors," *Softw. Eng. IEEE Trans.*, vol. 33, no. 1, pp. 2-13, 2007.
- [6] G. Chandrashekar, and F. Sahin, "A survey on feature selection methods," *Comput. Electr. Eng.*, vol. 40, no. 1, pp. 16-28, 2014.
- [7] مالکی، پیش‌بینی خطای نرم‌افزار بر اساس وابستگی‌های قطعات که با استفاده از داده‌کاوی، پایان‌نامه کارشناسی ارشد، دانشکده آموزش‌های الکترونیکی، دانشگاه شیراز، شیراز، ایران، ۱۳۹۰.
- [8] E. Alpaydin, *Introduction to Machine Learning*, MIT press, 2014.
- [9] S. Shivaji, S. Member, E. J. W. Jr, and S. Member, "Reducing features to improve code change-based bug prediction," *IEEE Trans. Softw. Eng.*, vol. 39, no. 4, pp. 552-569, 2013.

مشترک در هر دو ارزیابی A1 و A2 هستند. این ویژگی‌ها نیز در جدول ۶ آورده شده‌اند.

نتایج حاصل از روش انتخاب ویژگی پیشنهادی در جدول ۷، با نتایج پنج پژوهش دیگر مقایسه شده است. از هر یک از این پژوهش‌ها، نتایج روش‌هایی که بهترین عملکرد را داشته‌اند، بیان شده است. در این جدول بهترین نتایج برای هر یک از دادگان برجسته شده است. ما در این پژوهش از ۸ دادگان استفاده کرده‌ایم، اما با توجه به این‌که در اکثر پژوهش‌های مرتبط تنها از تعداد اندکی از دادگان استفاده شده است، برای برخی از دادگان توسط دیگر پژوهش‌ها نتایج اعلام نشده است، بنابراین خانه مربوطه برای این پژوهش‌ها در جدول با خط تیره مشخص شده است. نتایج مقایسه‌ها نشان‌دهنده این است که روش پیشنهادی این پژوهش در مقایسه با پژوهش‌های مشابه، در اکثر دادگان عملکرد بهتری در پیش‌بینی ماژول‌های نرم‌افزاری مستعد اشکال نسبت به پژوهش‌های مشابه داشته است.

جدول ۵: ویژگی‌های برگزیده توسط روش پیشنهادی

نماد ویژگی	نام کامل ویژگی	معادل فارسی
X1	Line count of code	تعداد کل خطوط کد
X9	Number of unique operators and operands	تعداد کل عملگرها و عملوندهای یکتا
X18	Number of comment-only lines	تعداد خطوط توضیحات از ماژول
X31	Number of parameter	تعداد پارامترهای ورودی
X36	Number of modified condition	تعداد شرط‌های اصلاح‌شده

جدول ۶: ویژگی‌های غیر مؤثر شناخته‌شده توسط روش پیشنهادی

نماد ویژگی	نام کامل ویژگی	معادل فارسی
X17	Number of blank lines	تعداد خطوط خالی ماژول
X29	Essential density	تراکم ضروری
X39	Normalized cyclomatic complexity	میزان نرمال‌شده تعداد مسیره‌های مستقل خطی در نمودار جریان ماژول

جدول ۷: مقایسه نتایج روش پیشنهادی با نتایج دیگر پژوهش‌ها

دادگان	PC3	PC2	PC1	MW1	MC2	MC1	KC3	CM1	این مقاله
A1	۰/۸۱	۰/۹۷	۰/۹۳	۰/۹۷	۰/۹۳	۰/۹۸	۰/۹۴	۰/۸۵	مقاله
A2	۰/۹۱	۰/۹۸	۰/۸۴	۰/۸۸	۰/۹۰	۰/۵۵	۰/۸۸	۰/۸۹	
[۱۶]	۰/۸۰	-	۰/۷۹	۰/۷۹	۰/۷۱	-	۰/۸۲	۰/۷۸	
[۲۲]	-	۰/۹۵	-	-	-	۰/۹۸	۰/۸۶	-	
[۲۶]	۰/۷۹	۰/۹۲	۰/۸۵	۰/۷۵	۰/۷۸	-	۰/۷۹	۰/۷۶	
[۳۲]	۰/۸۴	۰/۸۰	۰/۸۳	-	-	۰/۹۰	-	-	
[۳۳]	۰/۸۵	۰/۹۰	۰/۸۹	۰/۷۸	۰/۷۷	۰/۹۷	۰/۸۳	۰/۷۹	

۸- نتیجه‌گیری و پژوهش‌های پیش رو

در این پژوهش به‌دلیل اهمیت و کارایی روش‌های انتخاب ویژگی در بهبود عملکرد روش‌های یادگیری ماشین، روش‌های ترکیبی انتخاب

- [23] H. Wang, T. M. Khoshgoftaar, and A. Napolitano, "Software measurement data reduction using ensemble techniques," *Neurocomputing*, vol. 92, pp. 124-132, 2012.
- [24] H. Wang, "A comparative study of ensemble feature selection techniques for software defect prediction," *Mach. Learn. Appl. (ICMLA), Ninth Int. Conf. on. IEEE*, pp. 135-140, 2010.
- [25] T. Choeikiwong, and P. Vateekul, "Software defect prediction in imbalanced data sets using unbiased support vector machine," *Inf. Sci. Appl. Springer Berlin Heidelberg*, vol. 339, pp. 923-931, 2015.
- [26] R. S. Wahono, N. Suryana, and S. Ahmad, "Metaheuristic optimization based feature selection for software defect prediction," *J. Softw.*, vol. 9, no. 5, pp. 1324-1333, 2014.
- [27] S. García, J. Luengo, and F. Herrera, *Data Preprocessing in Data Mining*, Springer International Publishing, Switzerland, Intelligent Systems Reference Library, vol. 72, pp. 163-193, 2015.
- [28] D. Radjenović, M. Heričko, R. Torkar, and A. Živković, "Software fault prediction metrics: a systematic literature review," *Inf. Softw. Technol.*, vol. 55, no. 8, pp. 1397-1418, Aug 2013.
- [29] E. Arisholm, L. C. Briand, and E. B. Johannessen, "A systematic and comprehensive investigation of methods to build and evaluate fault prediction models," *J. Syst. Softw.*, vol. 83, no. 1, pp. 2-17, 2010.
- [30] D. Rodriguez, I. Herraiz, R. Harrison, J. Dolado, and J. C. Riquelme, "Preliminary comparison of techniques for dealing with imbalance in software defect prediction," *18th Int. Conf. Eval. Assess. Softw. Eng. ACM*, 2014.
- [۳۱] علیقارداشی، پیش‌بینی اشکال در ماژول‌های نرم‌افزاری با تلفیق روش‌های انتخاب ویژگی در مدل‌های توسعه‌یافته ماشین بردار پشتیبان، پایان‌نامه کارشناسی ارشد، دانشکده مهندسی برق و کامپیوتر، دانشگاه یزد، یزد، ایران، ۱۳۹۴.
- [32] B. Shuai, H. Li, M. Li, Q. Zhang, and C. Tang, "Software defect prediction using dynamic support vector machine," *Comput. Intell. Secur. (CIS), 9th Int. Conf. IEEE*, pp. 260-263, 2013.
- [33] Z. Sun, Q. Song, and X. Zhu, "Using coding-based ensemble learning to improve software defect prediction," *IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews*, vol. 42, no. 6, pp. 1806-1817, 2012.
- [10] S. Srivastava, N. Joshi, and M. Gaur, "A review paper on feature selection methodologies and their applications," *Int. J. Comput. Sci. Netw. Secur.*, vol. 14, no. 5, p. 78, 2014.
- [11] K. Muthukumar, A. Rallapalli, and N. L. B. Murthy, "Impact of feature selection techniques on bug prediction models," *Proc. 8th India Softw. Eng. Conf.*, pp. 120-129, 2015.
- [12] J. Apolloni, G. Leguizamón, and E. Alba, "Two hybrid wrapper-filter feature selection algorithms applied to high-dimensional microarray experiments," *Applied Soft Computing*, vol. 38, pp. 922-932, 2016.
- [13] A. Idris, and A. Khan, "Churn prediction system for telecom using filter-wrapper and ensemble classification," *The Computer Journal*, bvx123, 2016.
- [14] S. Huda, J. Abawajy, M. Alazab, M. Abdollahian, R. Islam, and J. Yearwood, "Hybrids of support vector machine wrapper and filter based framework for malware detection," *Future Generation Computer Systems*, vol. 55, pp. 376-390, 2016.
- [15] K. Gao, T. M. Khoshgoftaar, H. Wang, and N. Seliya, "Choosing software metrics for defect prediction: an investigation on feature selection techniques," *Softw. Pract. Exp.*, vol. 41, no. 5, pp. 579-606, 2011.
- [16] S. Liu, X. Chen, W. Liu, J. Chen, Q. Gu, and D. Chen, "FECAR: a feature selection framework for software defect prediction," *Comput. Softw. Appl. Conf. (COMPSAC), IEEE 38th Annu.*, pp. 426-435, 2014.
- [17] Z. Wang, Y. Shao, and T. Wu, "A GA-based model selection for smooth twin parametric-margin support vector machine," *Pattern Recognit.*, vol. 46, no. 8, pp. 2267-2277, 2013.
- [18] Z. Zhao, F. Morstatter, S. Sharma, S. Alelyani, A. Anand, and H. Liu, "Advancing feature selection research," *ASU Featur. Sel. Repos.*, pp. 1-28, 2010.
- [19] A. Jović, and N. Bogunović, "A review of feature selection methods with applications," *Info. Comm. Tech., Elec. Microelec. (MIPRO), 38th International Convention on*, 2015.
- [20] H. H. Hsu, C. W. Hsieh, and M. D. Lu, "Hybrid feature selection by combining filters and wrappers," *Expert Syst. Appl.*, vol. 38, no. 7, pp. 8144-8150, 2011.
- [21] J. Tang, S. Alelyani, and H. Liu, "Feature selection for classification: a review," *Data Classif. Algorithms Appl.*, p. 37, 2014.
- [22] I. H. Laradji, M. Alshayeb, and L. Ghouti, "Software defect prediction using ensemble learning on selected features," *Inf. Softw. Technol.*, vol. 58, pp. 388-402, 2015.

زیرنویس‌ها

¹² Area Under ROC (Receiver Operating Characteristic) Curve¹³ True Positive¹⁴ False Positive¹⁵ True Negative¹⁶ False Negative¹⁷ True Positive rate¹⁸ False Positive rate¹⁹ k-fold cross validation²⁰ Smooth Twin Parametric-Margin SVM¹ Fault-prone² Not fault-prone³ Forward selection⁴ Ensemble learning⁵ Greedy forward selection⁶ Meta-heuristic optimization⁷ Genetic algorithm⁸ Particle swarm optimization⁹ Bagging¹⁰ Symmetric uncertainty¹¹ Information gain