



## تشخیص کدهای مشکوک در کد منبع با استفاده از تکنیک های داده کاوی

مینا ریاحی عالم<sup>۱</sup>، رضا اکبری<sup>۲</sup>

<sup>۱</sup> کارشناس ارشد، مهندسی کامپیوتر- نرم افزار، دانشگاه صنعتی شیراز، شیراز،  
m.riahialam@gmail.com

<sup>۲</sup> استادیار، دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی شیراز، شیراز،  
akbari@sutech.ac.ir

### چکیده

با توجه به گسترش نرم افزارها، مهندسين نرم افزار و محققين به دنبال کاهش هزینه های تولید نرم افزار و افزایش کیفیت آن از جمله قابلیت توسعه، قابلیت پیمانه ای، استفاده مجدد، پیچیدگی، نگهداری و کارایی می باشند. بازسازی نرم افزار یکی از روش های موثر در افزایش کیفیت نرم افزار می باشد. یکی از نکات مورد نظر در بازسازی، تشخیص کدهای مشکوک در نرم افزار است. کد های مشکوک به علت ناکارآمد بودن طراحی یا پیاده سازی ناقص می باشند که تکامل، درک و حفظ نرم افزار را دشوار می نمایند. از این روش هایی جهت تشخیص کدهای مشکوک در نرم افزار مورد مطالعه قرار گرفته اند. یکی از روش های مورد استفاده، الگوریتم های داده کاوی می باشند که به بهبود صحت در تشخیص کدهای مشکوک کمک می نمایند. اما یکی از مسائلی که در داده کاوی مطرح است انتخاب ویژگی مناسب می باشد. روش های مختلفی از جمله الگوریتم های فراابتکاری جهت این امر مناسب می باشند. در این مقاله با استفاده از الگوریتم جستجوی فاخته، ابتدا عمل انتخاب ویژگی انجام شده است، سپس با استفاده از الگوریتم های داده کاوی، الگوریتم ها جهت تشخیص کدهای مشکوک بررسی شده اند.

### کلمات کلیدی

کد مشکوک، داده کاوی، انتخاب ویژگی، الگوریتم های فراابتکاری، الگوریتم جستجوی فاخته.

خواهد یافت. اما برخی مواقع به دلیل خطاهای موجود در مراحل توسعه ی نرم افزار ایراداتی ایجاد می شود که هزینه بر هستند. یکی از راه های افزایش کیفیت نرم افزار بر طرف نمودن عیوب نرم افزار در مراحل مختلف توسعه می باشد. اگر این عیوب برطرف نشوند، باعث به وجود آمدن بدهی فنی<sup>۱</sup> می شود [۱]. پیامد های بدهی فنی بر کیفیت نرم افزار به طور تجربی ثابت شده است. به عنوان مثال طراحی ضعیف یا پیاده سازی ضعیف یکی از موارد ایجاد بدهی ضعیف خواهد شد و امکان دارد بر روی نگهداری یک سیستم نرم افزاری اثر بگذارد.

### ۱- مقدمه

امروزه با افزایش وسعت نرم افزارها، فرآیند توسعه نرم افزار نیز گسترش یافته است. فرآیند توسعه ی نرم افزار می تواند شامل مستند سازی، تست نرم افزار، توسعه ی جدید، نگهداری و استفاده ی مجدد باشد. با توجه به پیشرفت دنیای تکنولوژی امروزه نیاز به نرم افزارهایی با کیفیت بیشتر می باشد. به همین جهت کیفیت نرم افزار با معیارهای گوناگونی مورد ارزیابی قرار می گیرد. همچنین راه های مختلفی جهت افزایش کیفیت نرم افزارها موجود می باشند. اگر مراحل مختلف توسعه ی نرم افزار به خوبی صورت گیرد، کیفیت افزایش



VFI و J48، صحت کمتری نسبت به الگوریتم‌های در نظر گرفته شده جهت پیش بینی کد مشکوک دارند.

یکی از روش‌های مورد استفاده در تشخیص کدهای مشکوک با استفاده از الگوریتم‌های داده کاوی، [۶] می‌باشد. در این روش شش الگوریتم یادگیری ماشین موجود در وکا شامل J48، SMO، جنگل تصادفی، JRip و شبکه‌ی بیزی بر روی مجموعه داده کوالیتاس کرایس<sup>۱۴</sup> به کار برده است. کوالیتاس کرایس یک مجموعه داده‌ی کاملی می‌باشد که شامل ۱۱۱ نرم افزار است. همچنین یکی از مزیت‌های دیگر این مجموعه داده این است که نسخه‌های مختلف هر یک از نرم افزارها را در نظر گرفته است. در این روش داده‌ها از ۷۴ سیستم نرم افزاری می‌باشند که از مجموعه داده‌ی کوالیتاس کرایس گرفته شده‌اند. طبق پیاده‌سازی انجام شده جهت تشخیص چهار کد مشکوک، خصیصه حسادت، کلاس داده، کلاس بزرگ و متد بزرگ، J48، SMO، جنگل تصادفی و JRip صحتی بیش از ۹۰ درصد بر روی کل داده‌ها و به طور متوسط بهترین کارایی را نسبت به الگوریتم شبکه‌ی بیزی داشته‌اند.

در [۷] نیز که در ادامه تحقیقات پیشین بر روی الگوریتم‌های داده کاوی در [۶] می‌باشد، ۱۶ الگوریتم مختلف یادگیری ماشین از جمله B-J48 هرس شده، B-J48 هرس نشده، B-J48 کاهش خطای هرس، B - درخت تصادفی، B-JRip، B بیز ساده، B-SMO RBF Kernel و B-SMO LibSVM C-SVC، LibSVM C-SVC کرنل خطی<sup>۱۵</sup>، LibSVM C-SVC کرنل شعاعی<sup>۱۶</sup>، LibSVM C-SVC کرنل ساگمویند<sup>۱۷</sup>، LibSVM v-SVC کرنل خطی، LibSVM v-SVC کرنل چند جمله‌ای، LibSVM v-SVC کرنل شعاعی و LibSVM v-SVC کرنل ساگمویند به کار رفته است که در این پژوهش می‌توان به اضافه شدن الگوریتم بوستینگ به الگوریتم‌های گذشته اشاره نمود که تأثیر به‌سزایی در افزایش صحت الگوریتم‌ها دارا بوده است. این پژوهش همانند کار قبل بر روی چهار کد مشکوک، خصیصه حسادت، کلاس داده، کلاس بزرگ و متد بزرگ به وسیله‌ی نرم افزار وکا اجرا و بر روی ۷۴ سیستم نرم افزاری در مجموعه داده کوالیتاس کرایس در نظر گرفته شده است که الگوریتم‌ها در اعتبارسنجی ضربدری به بالاترین کارایی دست یافته‌اند که بالاترین کارایی به وسیله J48 و جنگل تصادفی بدست آمده است و بدترین کارایی توسط SVM کسب شده است. همچنین میزان صحت در تشخیص کدهای مشکوک نیز در نظر گرفته شده است. یکی از ایرادات این روش عدم انتخاب ویژگی بر روی داده‌های بدست آمده از منبع کد نرم افزارها می‌باشد. یکی از مسائلی که در این روش در نظر گرفته نشده است، مسئله‌ی انتخاب ویژگی می‌باشد.

اما در برخی از تحقیقات صورت گرفته، از الگوریتم‌های بهینه‌سازی<sup>۱۸</sup> جهت انتخاب ویژگی استفاده شده است. یکی از الگوریتم‌های مورد بررسی الگوریتم جستجوی فاخته می‌باشد که در پژوهش‌های مختلف جهت انتخاب ویژگی به کار رفته است. در پژوهش [۸]، به استفاده از الگوریتم جستجوی فاخته برای انتخاب ویژگی بر روی چهار نوع مجموعه داده پرداخته است. در این روش با تغییر در تابع برازش الگوریتم به بهبودی آن در روند انتخاب ویژگی کمک نموده است.

در پژوهش‌های مختلف الگوریتم جستجوی فاخته بر روی مجموعه داده‌های مختلف به خصوص داده‌های پزشکی مورد استفاده قرار گرفته

از این رو جهت افزایش کیفیت نرم افزار نیاز است تا اطلاعاتی در ارتباط با نرم افزار کسب نمود. مخازن نرم افزار<sup>۲</sup> شامل اطلاعات زیادی پیرامون نرم افزار می‌باشند که با استخراج اطلاعات مفید می‌توان به روند تکامل نرم افزار کمک نمود. یکی از مخازن نرم افزار مورد استفاده منبع کد است [۲]. بدین صورت می‌توان با بدست آوردن متریک‌های مناسب از منبع کد، به رفع نواقص موجود در پیاده‌سازی کمک نمود که این کار باز سازی<sup>۳</sup> نرم افزار می‌باشد. در مفهوم تکامل نرم افزار<sup>۴</sup>، بازسازی نیز به بهبود کیفیت نرم افزار از جمله قابلیت توسعه<sup>۵</sup>، قابلیت پیمانه‌ای<sup>۶</sup>، استفاده مجدد، پیچیدگی، نگهداری و کارایی<sup>۷</sup>، به کار می‌روند [۳]. به این نکته هم اشاره شود که برخی مواقع نیز یک پیاده‌سازی ناقص ناشی از یک طراحی ضعیف می‌باشد.

بازسازی در نرم افزار به این صورت است که ابتدا بر روی منبع کد تحلیل صورت می‌گیرد. سپس بررسی می‌شود که آیا کد مورد نظر نیاز به بازسازی دارد یا خیر. در واقع در این قسمت کدهای مشکوک<sup>۸</sup> تشخیص داده می‌شوند. در پایان نیز اگر نیاز باشد، بازسازی صورت می‌پذیرد. نرم افزارها دارای معیارهای مختلفی می‌باشند که می‌توان با استفاده از آنها ویژگی‌های مختلف یک نرم افزار مشخص نمود.

همانطور که بیان شد، منبع کد دارای اطلاعاتی است که می‌توان به وسیله‌ی در نظر گرفتن متریک‌های مناسب به استخراج اطلاعات مفید در جهت بهبود کیفیت نرم افزار پرداخت. بدین منظور نرم افزارهایی برای محاسبه‌ی متریک‌ها از کد نرم افزار در دسترس می‌باشند. به عنوان نمونه می‌توان به نرم افزار اکلپس<sup>۹</sup> اشاره نمود. به منظور تشخیص کدهای مشکوک در نرم افزار متریک‌هایی از کد در نظر گرفته شده‌اند. از جمله متریک‌های مورد توجه در نرم افزار می‌توان به تعداد خطوط کد یک عملگر<sup>۱۰</sup>، تعداد خطوط کد، تعداد بسته<sup>۱۱</sup>، تعداد کل کلاس‌های در یک سیستم، در یک بسته و یا یک کلاس و غیره اشاره نمود [۴].

در ادامه ابتدا مروری بر تحقیقات پیشین که بر روی کدهای مشکوک و روش‌های مختلف جهت انتخاب ویژگی پرداخته می‌شود، سپس در بخش دوم، مفاهیم اولیه‌ی مورد نیاز در این پژوهش بیان شده است. بخش سوم نیز شامل روش پیشنهادی جهت تشخیص کدهای مشکوک و نحوه‌ی پیاده‌سازی و نتایج بدست آمده می‌باشد. در پایان نیز نتیجه‌گیری حاصل از این تحقیق صورت گرفته است.

## ۱-۱- مروری بر تحقیقات پیشین

با توجه به تأثیر کد مشکوک در نرم افزار، تحقیقات مختلفی پیرامون روش‌های تشخیص آنها انجام شده است. در [۵]، مانیرات<sup>۱۲</sup> و همکارانش به پیش‌بینی کدهای مشکوک از جمله کلاس تنبل، خصیصه حسادت، تابع طولانی، لیست پارامتر طولانی، دستور سوئیچ، زنجیر پیام و مرد واسطه می‌پردازد. به همین منظور از مجموعه داده‌هایی استفاده شده است که به وسیله‌ی هفت الگوریتم یادگیری ماشین بیز ساده، رگرسیون لجستیک<sup>۱۳</sup>، J48، VFI، IBk، IB1 و جنگل تصادفی به وسیله‌ی نرم افزار وکا آموزش و تست شده‌اند. داده‌های مورد استفاده با توجه به معیارهای مختلف نرم افزار مانند پیچیدگی، اطلاعات پایه، ارث بری و غیره در نظر گرفته شده‌اند. در نتایج بررسی شده در این تحقیق، مشخص شد که الگوریتم‌های بیز ساده،



کلاس داده دارای فیلدهایی است که متد های get و set برای این فیلدها می باشند و چیز دیگری ندارند. این قبیل کلاس های نگهدارنده های داده هستند و معمولاً توسط کلاس های دیگر مورد دستکاری قرار می گیرند. اگر در اوایل مراحل کار کلاس ها دارای فیلدهای عمومی باشند، می توان به وسیله ی به کارگیری سریع کپسوله سازی فیلد<sup>۳۳</sup> آن را صحیح نمود.

خصیصه حسادت زمانی رخ می دهد که یک متد به جای کلاس خود، بیشتر به کلاس دیگری تمایل دارد. اغلب تمرکز خصیصه حسادت بر روی داده می باشد. این قسمت را می توان با جابه جایی کد انجام داد.

تابع طولانی زمانی رخ می دهد که متدی طولانی باشد. برنامه هایی که دارای متد های کوتاه تری هستند عمر بهتر و بیشتری دارند. متدهایی که طولانی می باشند، سخت تر قابل فهم هستند. در زبان های برنامه نویسی قدیمی تر فراخوانی های زیرروال<sup>۳۴</sup> دارای سربار بودند. به همین دلیل برنامه نویسان از استفاده ی متدهای کوچک صرف نظر می کردند. اما زبان های شیء گرای جدید سربار فراخوانی در فرآیند را تا حدی حذف نموده اند.

کلاس بزرگ بدین صورت است که زمانی که یک کلاس سعی می کند کارهای زیادی انجام دهد، اغلب دارای متغیرهای زیادی می باشد. زمانی که تعداد متغیرها زیاد باشد، کدهای تکرار شونده اجتناب ناپذیرند و به همین ترتیب کلاس بزرگ می شود.

## ۲-۲- داده کاوی

با استفاده از الگوریتم های داده کاوی می توان ابتدا به آموزش نمونه های کد مشکوک موجود که برای تشخیص از منبع کد مورد نظر می باشند، استفاده نمود. از جمله ی این الگوریتم ها می توان به موارد زیر اشاره نمود.

• درخت تصمیم که یک الگوریتم با ناظر<sup>۳۵</sup> جهت دسته بندی می باشد. درخت تصمیم از ریشه<sup>۳۶</sup> آغاز می شود، سپس به شاخه<sup>۳۷</sup> ها که شامل گره<sup>۳۸</sup> ها می باشند، می شود و در نهایت تصمیم گیری در برگ<sup>۳۹</sup> صورت می گیرد[۱۴]. یکی از الگوریتم های مورد استفاده C4.5 می باشد که یک پیاده سازی کمی تغییر یافته از الگوریتم ID3 است. این الگوریتم با هرس درخت به بالا بردن صحت یادگیری کمک می نماید و باعث کاهش شاخه های تکراری می شود. از جمله الگوریتم های مورد استفاده می توان به J48 و جنگل تصادفی<sup>۴۰</sup> اشاره نمود.

• ماشین بردار پشتیبان<sup>۴۱</sup> یا SVM یکی از روش های یادگیری با نظارت است که از طبقه بندی و رگرسیون<sup>۴۲</sup> استفاده می نماید.

• JRip، یک یادگیرنده مبتنی بر استنتاج و مبتنی بر قواعد (RIPPER<sup>۴۳</sup>) است که تلاش می کند قوانین گزاره ای را که می تواند برای طبقه بندی ویژگی ها استفاده شود، ارائه می دهد[۱۵].

است. به عنوان نمونه در [۹]، از این الگوریتم جهت دسته بندی داده های سرطان استفاده کرده است. به دلیل اینکه بیان ژن<sup>۴۴</sup> دارای ویژگی های زیاد و نمونه های کم می باشد، با توجه به معیارهای ارزیابی محاسبه شده برای این الگوریتم، نشان می دهد که الگوریتم جستجوی فاخته مناسب برای طبقه بندی این نوع داده ها می باشد. همچنین در [۱۰]، الگوریتم جستجوی فاخته جهت طبقه بندی داده های سرطان سینه مورد استفاده قرار گرفته است.

• در بررسی های صورت گرفته بر روی داده های مختلف، الگوریتم فاخته عملکرد دارای عملکرد مناسبی می باشد. طبق مشاهدات این الگوریتم برای فضای جستجوی بزرگ و محیط های پویا و در حال رشد مناسب است.

## ۲- بیان مسئله

همانطور که بیان شد، با توجه به افزایش وسعت نرم افزارها، مهندسی نرم افزار و محققان به دنبال کاهش هزینه های تولید نرم افزار و افزایش کیفیت های مختلف نرم افزار می باشند. یکی از موارد مؤثر در کیفیت نرم افزار، منبع کد است که با توجه به افزایش وسعت نرم افزار، افزایش می یابد. به همین منظور مسئله ای که مطرح می شود این است که چگونه می توان روشی بهینه و کارآمد جهت بهبود منبع کد به کار بست.

کدهای مشکوک یکی از مشکلات مؤثر بر کیفیت نرم افزار هستند که به علت ناکارآمد بودن طراحی یا پیاده سازی ناقص می باشد[۸]. از این رو مسئله ی شناسایی کدهای مشکوک مهم می باشد. هر کد مشکوک به بررسی نوع خاصی از عناصر سیستم (به عنوان مثال، کلاس و یا متد) می پردازد که توسط ویژگی های آن بررسی شده است در این پژوهش، به بررسی روشی بهینه جهت تشخیص کدهای مشکوک از منبع کد پرداخته خواهد شد.

یکی از روش های تشخیص کدهای مشکوک، استفاده از الگوریتم های داده کاوی می باشد. یکی از چالش های مهم در داده کاوی و یادگیری ماشین، مسئله ی انتخاب ویژگی می باشد. در واقع انتخاب ویژگی فرآیندی است که در آن زیرمجموعه ای از داده های مرتبط، از داده های اصلی انتخاب می شوند[۱۱]. این کار سبب بهبود کارایی یادگیری، کاهش هزینه محاسبات و بهبود فهم مدل می شود.

## ۲-۱- کد مشکوک

در بازیابی نرم افزار یکی از نکات اصلی، تشخیص قسمت هایی از کد است که دارای ضعف در طراحی یا پیاده سازی است و نیازمند تغییرات می باشند. طبق تعریف مارتین فولر یک کد مشکوک به معنای یک طراحی ضعیف است[۱۲]. کد مشکوک مانند یک خطای نحوی<sup>۴۵</sup> یا اخطارهای کامپایلر<sup>۴۶</sup> نمی باشد. در واقع به صورت اعلام یک مشکل در طراحی یا برنامه نویسی به صورت عملی است که ممکن است در آینده مسئله ایجاد نماید. تشخیص کد مشکوک به صورت دستی در پروژه های نرم افزاری بزرگ، زمان بر می باشد. از این گذشته این کار هزینه بر و مستعد خطا است. زیرا به تجربه ی توسعه دهنده و میزان درک آن از برنامه وابسته است. در [۱۳]، ماتیتالا<sup>۴۷</sup> و همکارانش به دسته بندی کدهای مشکوک با توجه به اشتراکاتی که میان آن ها می باشد، پرداخته اند. از جمله کدهای مشکوک مورد بررسی در این مقاله می توان به کلاس داده، خصیصه ی حسادت، کلاس بزرگ و تابع طولانی اشاره نمود.

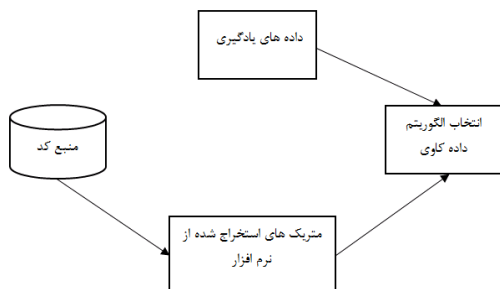
هایی که زیاد ایده آل نمی باشند، در لانه ها جایگزین شود. در ساده‌ترین حالت، هر لانه یک تخم دارد. الگوریتم این توانایی را دارد که برای مسائل پیچیده تر بسط داده شود. دلیل اینکه این الگوریتم عملکرد بهتری نسبت به برخی الگوریتم های دیگر دارد این است که دارای دو قابلیت جستجو می باشد. جستجوی محلی در  $1/4$  زمان جستجو و جستجوی عمومی  $3/4$  در زمان جستجو انجام می شود. این نشان دهنده ی آن است که در فضای عمومی بیشتر عمل کاوش صورت می پذیرد. این حالت زمانی رخ می دهد که  $P_a = 0.25$  در نظر گرفته شود.

## ۲-۴- تشخیص کد با استفاده از تکنیک های داده

### کاوی

به طور کلی روش تشخیص کد مشکوک با استفاده از الگوریتم های یادگیری ماشین از منبع کد را می توان به سه مرحله تقسیم نمود، شکل (۱):

- انتخاب داده های یادگیری مناسب که جهت آموزش به الگوریتم داده کاوی می باشند.
- آنالیز منبع کد و استخراج ویژگی های مناسب
- انتخاب الگوریتم مناسب با توجه به کد مشکوک مورد نظر.



شکل (۱): روش کلی تشخیص کد مشکوک با استفاده از الگوریتم های داده کاوی

## ۳- روش پیشنهادی تشخیص کد مشکوک

در روش های تشخیص مبتنی بر داده کاوی، جهت انتخاب، روشی مناسب تر است که بتوان بهترین زیر مجموعه را انتخاب و ایجاد نماید. از جمله راه هایی که می توان انتخاب ویژگی مناسبی از مجموعه ویژگی های داده انجام داد، روش های فراابتکاری می باشند. روش های فراابتکاری به دنبال یافتن بهینه ترین نقطه در فضای جستجو می باشد. از جمله الگوریتم مورد بررسی، الگوریتم جستجوی فاخته می باشد. الگوریتم جستجوی فاخته به دلیل سرعت، همگرایی مؤثر، پیچیدگی کمتر در پیاده سازی و پارامترهای ورودی یکی از الگوریتم های مناسب جهت انتخاب ویژگی استفاده نمود.

- بیز ساده<sup>۳۴</sup> یکی از روش های دسته بند کننده ی با ناظر است. بر مبنای تئوری بیز است. در این روش با در نظر گرفتن احتمال های شرطی موجود در بین ویژگی عمل می نماید.
- آدابوست<sup>۳۵</sup> از دسته ی الگوریتم های بوستینگ<sup>۳۶</sup> است [۱۶]. بوستینگ یک الگوریتم آموزشی است که جهت افزایش کارایی الگوریتم های یادگیری ماشین مورد استفاده قرار می گیرد. آدابوست برای داده های نامتوازن<sup>۳۷</sup> و پرت به کار گرفته می شوند و نتایج مربوط به این دسته از داده ها را بهبود می بخشند. این الگوریتم با کلاسه بندی داده ها به آموزش می پردازد.

## ۲-۳- الگوریتم های بهینه سازی

الگوریتم های بهینه سازی یا جستجوی محلی به دنبال یافتن پاسخی بهینه از مجموعه ی فضای جستجو می باشند [۱۷]. این الگوریتم ها به دو دسته ی الگوریتم های دقیق<sup>۳۸</sup> و الگوریتم های تقریبی<sup>۳۹</sup> تقسیم می شوند. اما الگوریتم های تقریبی می توانند جواب های خوب و نزدیک به بهینه را در زمان حل کوتاه بیابند. این الگوریتم ها برای مسائل بهینه سازی سخت، مناسب می باشند. یکی از الگوریتم های تقریبی، الگوریتم فراابتکاری می باشد [۱۸]. دو جزء اصلی در هر الگوریتم فراابتکاری، انتخاب بهترین راه حل و تصادفی بودن آن است. الگوریتم را می توان انتخاب بهترین جواب اشاره به نزدیک بودن راه حل ها به جواب مطلوب می باشد. همچنین تصادفی بودن آن نشان دهنده ی عدم به دام افتادن در نقطه ی بهینه محلی است. در این تحقیق نیز از این دسته از الگوریتم ها استفاده شده است. الگوریتم جستجوی فاخته از الگوریتم های فراابتکاری می باشد. الگوریتم جستجوی فاخته یا CS روش بهینه سازی فراابتکاری است که رویکردی تکاملی در جستجوی راه حل بهینه دارد. این الگوریتم در سال ۲۰۰۹ توسط زین-شی یانگ<sup>۴۰</sup> و سواش دب<sup>۴۱</sup>، الهام گرفته شده از رفتار پرند ای به نام فاخته می باشد، توسعه یافته است [۱۹]. این الگوریتم اخیرا نشان داده است که از الگوریتم های ژنتیک و این بهینه سازی جمعیت ذرات<sup>۴۲</sup> مؤثر تر عمل کرده است. پرندگان فاخته به این صورت عمل می کنند که تخم های خود را در لانه های عمومی قرار می دهند. در این بین ممکن است حتی تخم های دیگر را به دلیل افزایش احتمال تبدیل تخم های خود به جوجه، از بین ببرند. تعداد زیادی از گونه ها، در لانه های میزبان تخم گذاری می کنند.

به طور کلی الگوریتم جستجوی فاخته به صورت زیر عمل می نماید:

- هر فاخته یک تخم یک مرتبه می گذارد و آن را در یک لانه ی تصادفی رها می کند.
- بهترین لانه ها با بالاترین کیفیت تخم به نسل های بعدی منتقل خواهد شد.
- تعداد لانه های میزبان موجود ثابت می باشد و تخم گذاشته شده توسط فاخته با احتمال  $P_a \geq (0.1)$  توسط فاخته ی میزبان کشف می شود. در این مورد، پرند ای میزبان می تواند تخم را بیرون کند یا به سادگی لانه را ترک کند و یک لانه ی کاملاً جدید ایجاد نماید.

جهت حل مسئله با الگوریتم فاخته می توان هر تخم در یک لانه نمایانگر یک راه حل است و هر فاخته فقط می تواند یک تخم بگذارد. هدف آن است تا از راه حل های بهتر و جدید بالقوه استفاده شود تا به جای راه حل

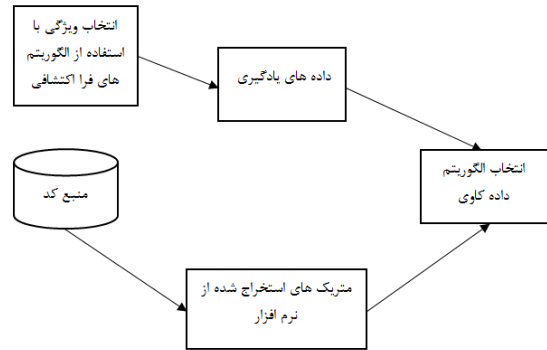
را دارا می باشند. در داده های تابع طولانی، B- جنگل تصادفی بهترین صحت را داراست.

بدین ترتیب الگوریتم B-جنگل تصادفی یکی از بهترین عملکردها را نسبت به سایر الگوریتم ها با انتخاب ویژگی توسط الگوریتم جستجوی فاخته داشته است.

#### ۴- نتیجه گیری

همانطور که بیان شد، مسئله ی بازسازی کد نرم افزار یکی از مواردی است که به افزایش کیفیت نرم افزار کمک می نماید. جهت بازسازی کد نرم افزار نیاز است تا ابتدا نقاطی از کد نرم افزار که نیاز به بازسازی دارند، تشخیص داده شوند. طبق بررسی های صورت گرفته در این پژوهش، با استفاده از الگوریتم جستجوی فاخته و تکنیک های مختلف داده کاوی مشخص شد که الگوریتم B-جنگل تصادفی با الگوریتم جستجوی فاخته دارای یکی از بهترین عملکرد نسبت به روش دیگر، در مسئله ی تشخیص کدهای مشکوک را داراست.

همچنین با توجه به بررسی ویژگی های انتخاب شده توسط الگوریتم جستجوی فاخته مشخص شد که این الگوریتم با انتخاب ویژگی های مربوط به همان کد مشکوک زیرمجموعه ای مناسب را انتخاب می نمایند.



شکل (۲): روش پیشنهادی تشخیص کد مشکوک با استفاده از الگوریتم های داده کاوی

#### ۳-۱- پیاده سازی روش پیشنهادی و نتایج

نحوه ی پیاده سازی به این صورت است که داده های مورد استفاده در این تحقیق با توجه به پژوهش [۷]، در نظر گرفته شده اند. داده ها از مجموعه ی کوالیتاس کراپس با استفاده نرم افزارهای استخراج متریک از نرم افزار ایجاد شده اند. همچنین این مجموعه داده با استفاده از سه نرم افزار مورد تحلیل قرار گرفته شده است تا کدهای مشکوک در آن ها تشخیص داده شود. سپس به صورت دستی داده ها چک شده اند که آیا تشخیص درست بوده است یا خیر. این مجموعه داده ها به چهار دسته ی کلاس داده، خصیصه حسادت، کلاس بزرگ و تابع طولانی تقسیم می شوند.

هر یک از داده ها مطابق و متریک های تعریف شده، دارای ویژگی هایی می باشند. دو مجموعه داده ی کلاس داده و کلاس بزرگ دارای ۶۶ ویژگی و ۴۲۰ نمونه و خصیصه حسادت و تابع طولانی دارای ۸۸ ویژگی و ۴۲۰ نمونه می باشند. روش مطالعه شده که مورد ارزیابی با روش پیشنهادی قرار خواهد گرفت، از ۱۶ الگوریتم یادگیری ماشین به تنهایی استفاده نموده است که نتایج نشان دهنده ی آن است که الگوریتم J48 و جنگلی تصادفی بالاترین میزان کارایی را دارا می باشند. این پژوهش بر روی الگوریتم هایی که صحت بالاتری نسبت به بقیه ی الگوریتم ها داشته اند که شامل الگوریتم B-J48 هرس شده، B-J48 هرس نشده، B-J48 کاهش خطای هرس، B- درخت تصادفی و J48 پیاده سازی شده اند. در این پژوهش نیز جهت کاهش ابعاد مجموعه ی داده از انتخاب ویژگی استفاده نموده است.

ابزار به کار رفته جهت پیاده سازی الگوریتم ها نرم افزار وکا نسخه ی ۳٫۸٫۱ است. ابتدا توسط الگوریتم جستجوی فاخته انتخاب ویژگی صورت می گیرد که تعداد تکرار جهت پیاده سازی جستجوی فاخته ۲۰ بار در نظر گرفته شده است. در جدول (۱)، (۲)، (۳) و (۴) به ترتیب نتایج صحت پیاده سازی روش پیشنهادی بر روی داده های کلاس داده، خصیصه ی حسادت، کلاس بزرگ و تابع طولانی که با الگوریتم جستجوی فاخته، ابتدا انتخاب ویژگی شده است، سپس با الگوریتم داده کاوی در نظر گرفته شده، نشان داده اند. در داده های کلاس داده، B- j48 هرس نشده و B-جنگل تصادفی بهترین صحت را نسبت به سایر الگوریتم ها با انتخاب ویژگی را دارا می باشند. در انتخاب ویژگی با الگوریتم جستجوی فاخته بر روی داده های خصیصه حسادت B-Jrip بهترین صحت را نسبت به سایر الگوریتم ها در همان انتخاب ویژگی را دارا می باشند. در انتخاب ویژگی با الگوریتم جستجوی فاخته B- جنگل تصادفی بهترین صحت را نسبت به سایر الگوریتم ها در همان انتخاب ویژگی

جدول (۱) صحت پیاده سازی با داده های کلاس داده

الگوریتم انتخاب ویژگی جستجوی فاخته	بدون انتخاب (مورد مطالعه)	نتایج	الگوریتم های داده کاوی
۹۸٫۰۴	۹۹٫۰۲	B- j 48 هرس شده	
۹۸٫۷۲	۹۸٫۶۷	B- j 48 هرس نشده	
۹۸٫۰۲	۹۸٫۰۷	B- j 48 کاهش خطای هرس	
۹۸٫۱	۹۸٫۸۳	B- Jrip	
۹۸٫۷۲	۹۸٫۵۷	B جنگل تصادفی	
۸۳٫۴۷	۹۷٫۳۳	B بیز ساده	
۹۰٫۰۵	۹۷٫۱۴	B- SMO RBF Kernel	
۹۴٫۹۳	۹۶٫۹	B- SMO کرنل های چند جمله ای	
۸۹٫۷۶	۹۸٫۵۵	j 48	
۹۷٫۰۲	۹۸٫۱۷	Jrip	





جدول (۳) صحت پیاده سازی با داده های تابع طولانی

بدون انتخاب(مورد مطالعه)	الگوریتم انتخاب ویژگی جستجوی فاخته	الگوریتم های داده کاوی	
۹۹,۴۳	۹۹,۵		B- j 48 هرس شده
۹۹,۲	۹۹,۶		B- j 48 هرس نشده
۹۹,۱۹	۹۹,۰۷		B- j 48 کاهش خطای هرس
۹۹,۰۳	۹۹,۱۶		B- Jrip
۹۹,۲۳	۹۹,۵۵		B جنگل تصادفی
۹۷,۸۶	۹۷,۱۶		B بیز ساده
۹۷	۹۷,۵۳		B- SMO RBF Kernel
۹۸,۶۷	۹۷,۲۹		B- SMO کرنل های چند جمله ای
۹۹,۱	۹۹,۲۷		j 48
۹۹,۰۲	۹۹,۱	Jrip	

جدول (۲) صحت پیاده سازی با داده های خصیصه حسادت

بدون انتخاب(مورد مطالعه)	الگوریتم انتخاب ویژگی جستجوی فاخته	الگوریتم های داده کاوی	
۹۶,۶۲	۹۶,۳۱		B- j 48 هرس شده
۹۶,۵	۹۶,۲۹		B- j 48 هرس نشده
۹۵,۹	۹۵,۷۴		B- j 48 کاهش خطای هرس
۹۶,۶۴	۹۶,۸		B- Jrip
۹۶,۴	۹۶,۷۹		B جنگل تصادفی
۹۱,۵	۸۴,۳		B بیز ساده
۹۳,۸۸	۸۹,۴۸		B- SMO RBF Kernel
۹۲,۰۵	۹۱,۰۵		B- SMO کرنل های چند جمله ای
۹۵,۹۵	۹۶		j 48
۹۵,۶۷	۹۵,۹۷	Jrip	

جدول (۳) صحت پیاده سازی با داده های کلاس بزرگ

بدون انتخاب(مورد مطالعه)	الگوریتم انتخاب ویژگی جستجوی فاخته	الگوریتم های داده کاوی	
۹۷,۰۲	۹۶,۹۸		B- j 48 هرس شده
۹۷,۰۲	۹۶,۹۵		B- j 48 هرس نشده
۹۷,۲۶	۹۷,۱۲		B- j 48 کاهش خطای هرس
۹۶,۹	۹۶,۸۸		B- Jrip
۹۶,۹۵	۹۷,۳۳		B جنگل تصادفی
۹۷,۵۴	۹۴,۲۱		B بیز ساده
۹۴,۶۲	۹۵,۱۴		B- SMO RBF Kernel
۹۴,۳۳	۹۶,۳۶		B- SMO کرنل های چند جمله ای
۹۷,۳۱	۹۶,۷۹		j 48
۹۷,۱۲	۹۷,۱۲	Jrip	

## مراجع

- [1] M. Tufano *et al.*, "When and Why Your Code Starts to Smell Bad," 2015, pp. 403–414.
- [2] A. E. Hassan, "The road ahead for Mining Software Repositories," 2008, pp. 48–57.
- [3] T. Mens and T. Tourwe, "A survey of software refactoring," *IEEE Trans. Softw. Eng.*, vol. 30, no. 2, pp. 126–139, Feb. 2004.
- [4] R. S. Pressman, *Software engineering: a practitioner's approach*, Eighth edition. New York, NY: McGraw-Hill Education, 2015.
- [5] N. Maneerat and P. Muenchaisri, "Bad-smell prediction from software design model using machine learning techniques," 2011, pp. 331–336.
- [6] F. A. Fontana, M. Zanoni, A. Marino, and M. V. Mantyla, "Code Smell Detection: Towards a Machine Learning-Based Approach," 2013, pp. 396–399.
- [7] C. C. Aggarwal, Ed., *Data classification: algorithms and applications*. Boca Raton: CRC Press, Taylor & Francis Group, 2014..
- [8] L. A. M. Pereira *et al.*, "A Binary Cuckoo Search and Its Application for Feature Selection," in *Cuckoo Search and Firefly Algorithm*, vol. 516, X.-S. Yang, Ed. Cham: Springer International Publishing, 2014, pp. 141–154.
- [9] C. Gunavathi and K. Premalatha, "Cuckoo search optimisation for feature selection in cancer classification:



20 Syntax  
21 Compiler  
22 Mäntylä  
23 Encapsulation  
24 Procedure  
25 Supervised  
26 Root  
27 Branch  
28 Node  
29 Leaf  
30 Random Forest  
31 Sector Vector Machine  
32 Regression  
33 Repeated incremental pruning to produce error  
reduction  
34 Naive Bayes  
35 Adaboost  
36 Boosting  
37 Imbalance  
38 Exact  
39 Approximate  
40 Xin-She Yang  
41 Suash Deb  
42 Particle swarm optimization  
43 Global Search

- a new approach,” *Int. J. Data Min. Bioinforma.*, vol. 13, no. 3, p. 248, 2015.
- [10] M. N. Sudha and S. Selvarajan, “Feature Selection Based on Enhanced Cuckoo Search for Breast Cancer Classification in Mammogram Image,” *Circuits Syst.*, vol. 07, no. 04, pp. 327–338, 2016.
- [11] Huan Liu and Lei Yu, “Toward integrating feature selection algorithms for classification and clustering,” *IEEE Trans. Knowl. Data Eng.*, vol. 17, no. 4, pp. 491–502, Apr. 2005.
- [12] M. Fowler and K. Beck, *Refactoring: improving the design of existing code*. Reading, MA: Addison-Wesley, 1999.
- [13] M. Mantyla, J. Vanhanen, and C. Lassenius, “A taxonomy and an initial empirical study of bad smells in code,” 2003, pp. 381–384.
- [14] Y. Zhao and Y. Zhang, “Comparison of decision tree methods for finding active objects,” *Adv. Space Res.*, vol. 41, no. 12, pp. 1955–1959, Jan. 2008.
- [15] [12] V. Rajeswari and K. Arunesh, “Analysing Soil Data using Data Mining Classification Techniques,” *Indian J. Sci. Technol.*, vol. 9, no. 19, May 2016.
- [16] M. Galar, A. Fernandez, E. Barrenechea, H. Bustince, and F. Herrera, “A Review on Ensembles for the Class Imbalance Problem: Bagging-, Boosting-, and Hybrid-Based Approaches,” *IEEE Trans. Syst. Man Cybern. Part C Appl. Rev.*, vol. 42, no. 4, pp. 463–484, Jul. 2012.
- [17] H. H. Hoos and T. Stützle, *Stochastic local search: foundations and applications*. San Francisco, CA: Morgan Kaufmann Publishers, 2005.
- [18] X.-S. Yang, *Engineering optimization an introduction with metaheuristic applications*. Hoboken, N.J: John Wiley, 2010.
- [19] I. Fister, X.-S. Yang, D. Fister, and I. Fister, “Cuckoo Search: A Brief Literature Review,” in *Cuckoo Search and Firefly Algorithm*, vol. 516, X.-S. Yang, Ed. Cham: Springer International Publishing, 2014, pp. 49–62.

## پانویس ها

- 1 Technical Debt  
2 Software Repositories  
3 Refactoring  
4 Software Evolution  
5 Extensibility  
6 Modularity  
7 Efficiency  
8 Code Smell  
9 Eclipse  
10 Operator  
11 Package  
12 Maneerat  
13 Logistic regression  
14 Qualitas Corpus  
15 Linear Kernel  
16 Ratial Kernel  
17 Sigmoid Kernel  
18 Optimization Algorithm  
19 Gene Expression