

ارائه یک روش یادگیری ویژگی ترکیبی مبتنی بر الگوریتم شبیه‌سازی تبرید و برنامه‌نویسی ژنتیک

(مطالعه موردی: تشخیص بدخیمی سرطان سینه)

رسول صادقی^۱، دانشجوی کارشناسی ارشد؛ فردین ابدالی محمدی^۲، استادیار

۱- گروه مهندسی کامپیوتر و فناوری اطلاعات، دانشکده فنی و مهندسی، دانشگاه رازی، کرمانشاه، ایران، rasool.sadeghi@stu.razi.ac.ir

۲- گروه مهندسی کامپیوتر و فناوری اطلاعات، دانشکده فنی و مهندسی، دانشگاه رازی، کرمانشاه، ایران، fardin.abdali@razi.ac.ir

چکیده: امروزه استفاده از ابزارهای یادگیری ماشین در حوزه‌های مختلف از جمله تشخیص بیماری‌ها در حال گسترش است. علت این امر را می‌توان عملکرد متغیر و متمایل به خطای انسان در مقابل عملکرد ثابت ابزارهای یادگیری ماشین در زمینه تشخیص و طبقه‌بندی دانست. حیاتی بودن تشخیص در حوزه‌هایی مانند پزشکی، نیاز به بهبود تشخیص با روش‌های یادگیری ماشین را توجیه می‌کند. از جمله روش‌های افزایش دقت در این زمینه، الگوریتم‌های کاهش ویژگی و یادگیری ویژگی هستند. در این مقاله با ارائه یک روش یادگیری ویژگی، دقت روش‌های مبتنی بر یادگیری ماشین افزایش یافته است. روش پیشنهادی شامل سه فاز افزایش کیفیت داده، انتخاب ویژگی و یادگیری ویژگی است. در فاز اول، مقادیر ازدست‌رفته با شاخص پراکندگی میانگین و یا مد جایگزین می‌شوند در فاز دوم، الگوریتمی مبتنی بر الگوریتم شبیه‌سازی تبرید برای کاهش ویژگی و یافتن بهترین زیرمجموعه از ویژگی‌ها ارائه شده است. در فاز نهایی نیز الگوریتمی مبتنی بر الگوریتم برنامه‌نویسی ژنتیک به منظور یادگیری ویژگی‌های متمایزکننده ترکیبی ارائه شده است. روش پیشنهادی با استفاده از دو مجموعه داده استاندارد WBCD و WDBC ارزیابی شده است. نتایج به دست آمده با آخرین دستاوردها مقایسه شده است که حاکی از عملکرد بهبودیافته الگوریتم پیشنهادی است.

واژه‌های کلیدی: یادگیری ماشین، یادگیری ویژگی، برنامه‌نویسی ژنتیک، کاهش ویژگی، سرطان سینه

A Combined Feature-Learning Method Based on Simulated Annealing Algorithm and Genetic Programming (Case Study: Malignant Breast Cancer Diagnosis)

R. Sadeghi¹, MSc student; F. Abdali Mohammadi², Assistant Professor

¹ Department of Computer Engineering and Information Technology, Faculty of Engineering, Razi University, Kermanshah, Iran, Email: rasool.sadeghi@stu.razi.ac.ir

² Department of Computer Engineering and Information Technology, Faculty of Engineering, Razi University, Kermanshah, Iran, Email: fardin.abdali@razi.ac.ir

Abstract: Nowadays using machine learning tools in different areas such as disease detection is expanding. Origins of this expansion can be found in humans' unstable performance and stable performance of machine learning tools. Criticality of detection in areas such as medical proves the need for improvement in machine learning methods. feature reduction and feature learning are two ways that cause to precision increment. In this paper precision of machine learning algorithms is increased by feature learning. The proposed method contains three steps: data quality increment, feature selection, and feature learning. In the first step missing values are replaced with mean or mode (distribution index). In the second step a simulated annealing-based algorithm is presented to utilized as feature selection process and finding the best subset of features. In the final step, a genetic programming algorithms is presented to do the feature learning step. The proposed method is evaluated on two benchmark datasets (WBCD and WDBC). The results show performance improvement in machine learning algorithms in terms of precision if the proposed method used.

Keywords: Machine learning, feature learning, genetic programming, feature reduction, breast cancer

تاریخ ارسال مقاله: ۱۳۹۵/۰۸/۱۹

تاریخ اصلاح مقاله: ۱۳۹۵/۱۰/۱۱

تاریخ پذیرش مقاله: ۱۳۹۵/۱۲/۰۹

نام نویسنده مسئول: فردین ابدالی محمدی

نشانی نویسنده مسئول: ایران - کرمانشاه - باغ ابریشم - دانشگاه رازی - دانشکده فنی و مهندسی.

۱- مقدمه

دسته بزرگ و مهمی از مسائلی که می‌تواند به کمک کامپیوتر پاسخ داده شود، مسائل مربوط به تشخیص بیماری‌ها است. مزایای استفاده از کامپیوتر در حل چنین مسائلی را می‌توان این‌گونه برشمرد: عملکرد و دقت ثابت و بدون تأثیرپذیری از محیط، سرعت ثابت و بالا در تشخیص، عدم محدودیت جغرافیایی و عدم محدودیت زمانی در استفاده. البته لازمه این امر، ارتقاء دقت الگوریتم‌های تشخیصی است. با توجه به آمار منتشرشده از منابع معتبر [۱]، سرطان سینه یکی از بیش‌ترین موارد تشخیص داده‌شده سرطان در ایران و سایر کشورها است. برای مثال بر اساس آمار صادره از موسسه تحقیقات، درمان و آموزش سرطان در ایران [۱] در سال ۱۳۸۷، سرطان سینه با ۸۶۱۶ مورد جدید، پس از سرطان پوست، دومین نوع عمده از سرطان در ایران تشخیص داده شده است. این آمار در کشور آمریکا [۲]، در سال ۲۰۱۴ با ۲۹۵۲۴۰ مورد، بیش‌ترین بدخیمی تشخیص داده‌شده در این کشور بوده است. از آن‌جایی که یکی از بیش‌ترین بدخیمی‌های تشخیص داده‌شده در ایران و سایر کشورها سرطان سینه است؛ مبارزه با این بیماری یک تلاش همگانی و یکپارچه را می‌طلبد. یکی از راه‌های کاهش آمار مرگ‌ومیر ناشی از این بیماری اجرای طرح‌های غربالگری عمومی است که نقش به‌سزایی در تشخیص به‌موقع این بیماری و به‌تبع آن نجات بیماران دارد. با توجه به آمار بالای ابتلای به این بیماری، اجرای چنین طرح عظیمی به‌صورت دستی و بدون دخالت ابزارهای مکانیزه بسیار زمان‌بر و مستعد خطاست. به‌طور قطع، ابزارهای یادگیری ماشین می‌توانند تأثیر به‌سزایی در افزایش سرعت و دقت و کاهش خطا در تشخیص این بیماری داشته باشند؛ همچنین با توجه به ماهیت الگوریتمی روش‌های پیشنهادی در زمینه تشخیص بیماری، روش‌های مذکور بر روی سایر بیماری‌ها نیز قابل تعمیم هستند.

ویژگی‌های مورد استفاده در مقاله پیش رو شامل مقادیری عددی است؛ که توسط متخصصین پزشکی از تصاویر ماموگرافی استخراج شده است. از طرفی عملکرد الگوریتم‌های یادگیری ماشین به‌شدت وابسته به ویژگی‌هایی است که یادگیری براساس آن‌ها انجام می‌شود [۳]؛ بنابراین استفاده از ویژگی‌های مناسب، لازمه استفاده بهینه از الگوریتم‌های یادگیری ماشین است. راهکارهایی که افزایش دقت الگوریتم‌های یادگیری ماشین را به دنبال دارد شامل کاهش ویژگی، نگاشت ویژگی، برطرف کردن همبستگی‌ها و وزن‌دار کردن ویژگی‌ها است. یکی از روش‌هایی که تمامی موارد گفته‌شده را به‌صورت ضمنی انجام می‌دهد، یادگیری ویژگی است. یادگیری ویژگی یکی از مباحث مطرح شده در حوزه یادگیری ماشین است. هدف از یادگیری ویژگی یافتن و ایجاد مجموعه‌ای از ویژگی‌ها است که منجر به اجرای بهینه و حتی سریع‌تر الگوریتم‌های یادگیری ماشین در فاز عملی خواهد شد. به بیان ساده می‌توان یادگیری ویژگی را یک تابع انتقال به‌صورت $T:F \rightarrow NF$ بیان کرد که در آن F بردار ویژگی‌های اولیه با مقادیر خام و NF بردار ویژگی‌هایی است که بر اساس F به دست آمده است.

برنامه‌نویسی ژنتیک (GP) در واقع شکل توسعه‌یافته‌ای از روش‌های یادگیری ماشین و الگوریتم ژنتیک (GA) است که می‌تواند در یادگیری ویژگی استفاده شود. این الگوریتم دارای وجوه تمایزی با GA نیز هست؛ از جمله اینکه برخلاف GA، در GP نیازی به مشخص بودن دامنه‌ی پاسخ در یک محدوده خاص نیست (در GA تعداد کاراکترهای رشته نهایی از قبل مشخص است). موضوعی که سبب انعطاف بیش‌تر این روش نسبت به GA شده این است که برخلاف GA که خروجی آن یک رشته باینری است؛ خروجی GP یک رابطه ریاضی بهینه است که بیش‌ترین برازندگی را با مسئله داشته و با کم‌ترین خطای ممکن، ورودی‌های مسئله را به خروجی‌های آن نگاشت می‌کند [۴]. این ویژگی GP باعث می‌شود که خروجی آن به راحتی در نرم‌افزارها و یا حتی سخت‌افزارهای بلادرنگ قابل پیاده‌سازی باشد. در این مقاله تلاش شده است که به کمک الگوریتم GP، یادگیری ویژگی صورت گرفته و دقت تشخیص و طبقه‌بندی تا حد مطلوب افزایش یابد.

در این مقاله به‌منظور کاهش ویژگی‌ها نیز تلاش شده ابتدا به کمک الگوریتم‌های بهینه‌سازی، ابعاد مسئله کاهش یابد. چرایی این امر را می‌توان زمان‌بر بودن فاز تغییر فضای مسئله به کمک الگوریتم GP بر اساس زیر مجموعه‌ای از ویژگی‌ها دانست؛ چراکه در فاز مذکور، ساخت ژن بر اساس تمامی ترکیب‌های ممکن از ویژگی‌ها بررسی می‌شود. لذا ناگزیر به کنار گذاشتن ویژگی‌های بی‌ارزش پیش از فاز مذکور هستیم و ابزارهای بهینه‌سازی و به‌طور خاص الگوریتم بهینه‌سازی شبیه‌سازی تبرید، گزینه مناسبی برای این منظور خواهد بود. مزیت مذکور علاوه بر تأثیر مثبتی است که کاهش ویژگی در کاربردهای عملی و کاهش ورودی‌های مورد نیاز مسئله و به‌تبع آن کاهش هزینه‌ها دارد. کاهش ویژگی به‌گونه‌ای انجام می‌شود که به‌جای استفاده از n کمیت اندازه‌گیری شده برای طبقه‌بندی، از k کمیت استفاده شود به‌گونه‌ای که $k < n$ باشد. این کار به تنهایی نقش به‌سزایی در کاهش هزینه‌های محاسباتی نیز خواهد داشت؛ چراکه حجم نمونه‌برداری و آزمایش‌ها را به k/n کاهش می‌دهد. ضمن اینکه سرعت گردش کار را در فاز عملی نیز به طرز چشمگیری افزایش می‌دهد. البته لازمه آنچه گفته شد، تضمین دقت در کنار کاهش ابعاد مسئله است.

نوآوری‌های این پژوهش به‌شرح زیر است:

- ۱- کاهش ویژگی به کمک الگوریتم شبیه‌سازی تبرید.
- ۲- یادگیری ویژگی به کمک الگوریتم برنامه‌نویسی ژنتیک، که ذاتاً یک الگوریتم برازش منحنی است.
- ۳- بهبود عملکرد الگوریتم برنامه‌نویسی ژنتیک با تغییر تابع برازش به تابعی بر اساس ماشین بردار پشتیبان خطی.
- ۴- استفاده از معیارهای ارزیابی جامع برای ارزیابی روش ارائه‌شده.

در ادامه، ابتدا در بخش دوم سوابق پژوهشی مورد بررسی قرار می‌گیرد. در بخش سوم توضیحاتی در مورد الگوریتم شبیه‌سازی تبرید، GP و پیش‌پردازش داده‌ها ارائه می‌شود. در بخش چهارم، روش

مورد استفاده ارزیابی شده و در انتها نتایج به دست آمده تجزیه و تحلیل شده و نتیجه گیری ها ارائه می شود.

۲- سابقه پژوهش

تلاش برای کاهش اثر بیش برآزش با کمک یادگیری ویژگی که در [۳] انجام شده و همچنین استفاده از یادگیری ویژگی در بازشناسی افراد بر اساس تصاویر گرفته شده از زوایای مختلف در [۵] نیز توانسته تأثیر مثبت یادگیری ویژگی را در بهبود الگوریتم های یادگیری ماشین به اثبات برساند. عملکرد مناسب الگوریتم های بهینه سازی در حذف ویژگی های کم ارزش نیز در [۶-۱۱] به اثبات رسیده است.

الگوریتم های یادگیری ماشین کارایی خود را در حل مسائل پزشکی ثابت کرده اند [۱۲-۱۳]. برای نمونه در [۱۴-۱۵] تلاش شده تا به ترتیب برای ایجاد یک طبقه بند مولکولی و پیش بینی پاسخ به درمان ضد سرطان از الگوریتم GA استفاده شود. نتایج این پژوهش ها کارایی این الگوریتم را در مسائل مشابه به وضوح تأیید می کند. الگوریتم های درخت تصمیم و شبکه عصبی نیز مزایای خود را نسبت به روش آماری رگرسیون در [۱۶] نشان داده اند. در [۱۷] نیز GA در کاهش ویژگی ها خروجی بهتری نسبت به آزمون های آماری داشته است. همچنین در [۱۸] الگوریتم درخت تصمیم عملکردی بهتر را نسبت به SVM نشان داده است.

تشخیص سرطانی بودن یا نبودن بافت ها بر اساس بیان ژنی کاری بود که در [۱۹] ارائه شد. این روش توانسته بافت سرطانی و سالم را با دقت ۹۳٪ تفکیک می کند. کشف رابطه بین بیان ژنی و سرطان های تخمدان، پروستات و ریه در [۲۰] انجام شده و روش پیشنهادی در آن توانسته با دقتی بین ۹۴٪ تا ۹۸٪ عمل طبقه بندی را انجام دهد. همچنین در [۲۱] تلاشی برای ایجاد قواعد تشخیص بیماری سرطان ریه صورت گرفته و سه الگوریتم Association Rules، Rough Set و Genetic Algorithm بر روی داده ها اجرا شده است. طبق نتایج به دست آمده الگوریتم های GA، AR و RS توانسته اند به گونه ای مشابه با قواعد تشخیص پزشکی عمل کنند. مشکل این روش ها نیاز به ابزارهای پیشرفته ژنتیکی و مولکولی است که پروتئین های بیان شده را اندازه گیری کنند.

یکی از راه های پیش بینی ابتلا به سرطان سینه استفاده از ابزاری تحت عنوان ترموگرام است. ترموگرام یک ابزار مانیتورینگ است و قادر است ابتلا به بیماری سرطان سینه را تا ۱۰ سال پیش از وقوع آن پیش بینی کند. داده های مربوط به این ابزار نیازمند تفسیر است و اگر تحلیلگر به درستی داده ها را تفسیر نکند ارزش و توانایی این ابزار زیر سؤال می رود. به منظور مکانیزه کردن تفسیر داده های ترموگرام و برای کاهش خطای انسانی در تفسیر اطلاعات آن، طرحی تحقیقاتی و عملی در [۲۲] صورت گرفته است. در این کار عده ای از محققان سنگاپوری از الگوریتم های داده کاوی برای استخراج قواعد تصمیم گیری استفاده کردند. تا به کمک آن بتوان نتایج ترموگرام را تفسیر کرد. آن ها از CLFNN (Complementary Learning Fuzzy Neural Network)

استفاده کردند. CLFNN الگوریتمی است قدرتمند و یکسری قواعد تصمیم گیری در اختیار می گذارد که به کمک آن به راحتی می توان داده های ابزار ترموگرام را با کم ترین خطا تحلیل کرد. در [۲۳-۲۴] نیز تلاش هایی برای تشخیص خوش خیمی سرطان سینه به کمک مجموعه داده WBCD صورت گرفته که روش های پیشنهادی توانسته عملکردی با دقت ۹۵/۵۶٪ تا ۹۵/۶۰٪ داشته باشد. در [۲۵-۲۶] نیز بیماری سرطان سینه به کمک روش های یادگیری ماشین و بر اساس داده های استخراج شده از تصاویر MRI با دقت ۹۶٪ تا ۹۸٪ تشخیص داده شده است. تلاش برای پیش بینی بقای بیماران مبتلا به این سرطان نیز در [۲۷]، دقتی برابر با ۸۳٪ را به دنبال داشته است. عملکرد حاصل از یک روش پیشنهادی در [۲۸] بر روی مجموعه داده WDBC شامل ۳۰ ویژگی مربوط به سرطان سینه نیز توانسته عملکردی برابر با ۹۷/۳۸٪ را به دنبال داشته باشد. پژوهش دیگری که بر روی همین مجموعه داده صورت گرفته است، ترکیب الگوریتم های ژنتیک، ازدحام ذرات و کلونی مورچه ها با الگوریتم ماشین بردار پشتیبان است [۲۹].

۳- روش پیشنهادی

روند کلی روش پیشنهادی به این صورت است که ابتدا کیفیت داده ها، با جایگزینی مقادیر از دست رفته، افزایش می یابد؛ پس از آن با بهره گرفتن از الگوریتم بهینه سازی شبیه سازی تیرید، زیرمجموعه بهینه از ویژگی ها به گونه ای انتخاب می شود که رابطه ۱ برقرار باشد.

$$SF = SA_{DATA}, \quad |SF| < |DATA| \quad (1)$$

در این رابطه SF ماتریس مقادیر زیرمجموعه ویژگی ها، SA تابع انتخاب ویژگی و DATA ماتریس داده های اولیه هستند. در گام بعدی از روش پیشنهادی، به کمک GP، ویژگی های جدیدی بر اساس ویژگی های انتخاب شده ساخته می شود. این گام به صورتی انجام می شود که در رابطه ۲ قابل ملاحظه است.

$$ND = GP_{SF} \quad (2)$$

در این رابطه ND بیانگر مجموعه ویژگی های جدیدی است که به کمک GP و بر اساس SF ایجاد شده است. GP نیز بیانگر تابع نگاشت فضای داده ها به کمک برنامه نویسی ژنتیک است. در واقع الگوریتم GP ژن هایی تولید می کند که هر ژن بیانگر یک رابطه ریاضی بر روی یک سری از ویژگی های ورودی بوده و هر ژن می تواند با اعمال روابط ریاضی بر روی تعدادی از ویژگی ها یک ویژگی جدید ایجاد کند و در صورت ایجاد چند ژن به کمک GP می توان چند ویژگی جدید ایجاد کرد.

روند کلی آنچه گفته شد به صورتی است که در الگوریتم ۱ دیده می شود. خطوط اول و دوم از الگوریتم ۱ به ترتیب بهبود کیفیت داده و کاهش ابعاد مسئله را انجام می دهند. همان طور که از خط ۱۰ الگوریتم برمی آید، شرط خاتمه الگوریتم، بررسی تمام زیرمجموعه های ممکن از ویژگی های کاهش یافته و یا رسیدن به حدی مطلوب از برآزش است. در خط ۶ از الگوریتم ۱ به کمک GP، فضای مسئله تغییر می یابد. چراکه

را مورد بررسی قرار می‌دهند، گزینه مناسبی هستند. یکی از این الگوریتم‌ها الگوریتم بهینه‌سازی شبیه‌سازی تبرید است. این الگوریتم از مباحث متالورژی الهام گرفته شده است. روشی که از آن الهام گرفته شده اساساً برای سرد کردن تدریجی فلزات به منظور رسیدن به بهترین آرایش و کم‌ترین انرژی است. در شبیه‌سازی تبرید هر نقطه از فضای جستجو به صورت یک حالت از سیستم ترمودینامیکی در نظر گرفته می‌شود. هر حالت نیز ممکن است به نقاط همسایه خود تغییر حالت دهد. این تغییر حالت اگر در راستای کاهش انرژی باشد به طور قطع انجام خواهد شد؛ اما اگر در راستای افزایش انرژی باشد آنگاه این تغییر حالت به صورت احتمالی انجام خواهد شد و این احتمال به دو عامل بستگی دارد: عامل اول اختلاف انرژی نقطه جدید با نقطه فعلی و عامل دوم دمایی است که تغییر حالت در آن صورت می‌گیرد. هر چقدر حالت جدید، انرژی بیشتری داشته باشد، احتمال تغییر به آن کاهش می‌یابد و هر چه دما در آن لحظه بالا باشد نیز این احتمال افزایش می‌یابد. رابطه ریاضی به کاررفته برای محاسبه احتمال بر اساس این دو پارامتر، می‌تواند تابع توزیع احتمال بولتزمن باشد (رابطه ۳).

$$e^{\left(\frac{-\Delta E}{T}\right)} \quad (3)$$

در این رابطه ΔE اختلاف انرژی یا همان اختلاف هزینه است که از طریق رابطه ۴ محاسبه می‌شود. پارامتر T نیز بیانگر دما است.

$$\Delta E = C_n - C_c \quad (4)$$

در این رابطه C_c همان هزینه پاسخ فعلی و C_n هزینه یک پاسخ جدید در الگوریتم شبیه‌سازی تبرید هستند. در الگوریتم شبیه‌سازی تبرید، که روند کلی آن در الگوریتم ۲ آمده است، در هر دما (گام) تعدادی نقاط همسایگی در اطراف پاسخ فعلی به صورت تصادفی ایجاد شده و پاسخ فعلی در راستای کاهش انرژی در بین این نقاط تغییر حالت می‌دهد. بررسی همسایگی‌ها در هر دما به تعداد مشخصی انجام شده و پس از این کار، دما کاهش می‌یابد و بررسی و یافتن همسایگی با انرژی کم‌تر برای دمای جدید انجام می‌شود. کاهش دما طبق رابطه ۵ انجام می‌شود.

$$T_{i+1} = \alpha \times T_i, \quad 0 < \alpha < 1, T_0 \neq 0 \quad (5)$$

نکته مهم، این است که در طول اجرای الگوریتم و با کاهش دما، به مرور امکان تغییر به حالات با انرژی بیشتر، کاهش می‌یابد؛ درحالی که این احتمال در ابتدای الگوریتم بالا است. این امر سبب می‌شود که در ابتدای الگوریتم، خاصیت جستجو یا (Exploration) بالا بوده و به مرور این خاصیت کاهش یافته و خاصیت بهره‌برداری (Exploitation) افزایش یابد و این سبب همگرایی تدریجی به سمت نقطه بهینه خواهد شد. نمایی کلی از روش کار الگوریتم شبیه‌سازی تبرید در شکل ۱ قابل مشاهده است. برای کاهش ویژگی به کمک الگوریتم شبیه‌سازی تبرید، یک زیرمجموعه k عضوی از بین ویژگی‌های موجود به گونه‌ای انتخاب می‌شود که کم‌ترین میزان طبقه‌بندی غلط را در اجرای درخت تصمیم بر روی

داریم: $GP: SF \rightarrow ND$. در این میان برازش‌های به دست آمده به ازای هر زیرمجموعه از ویژگی‌های ورودی به GP در بردار fit به طول n نگه‌داری می‌شود. این کار تا محقق شدن شرط خاتمه ادامه می‌یابد. در پایان، بهترین مجموعه داده ایجاد شده، با بیش‌ترین برازش، به عنوان داده‌ی نهایی تعیین شده و مدل بر اساس آن ایجاد می‌شود. این مجموعه داده در یکی از عناصر $D_{1..n}$ قرار داشته که اندیس آن از طریق خط ۱۱ از الگوریتم استخراج می‌شود.

الگوریتم ۱: روند کلی روش پیشنهادی

01. *Check data quality*
02. $Data_{n \times k} \leftarrow SAFR(Data_{n \times m}), \quad k < m$
03. $fitness \leftarrow 0$
04. $fit_{1..n} \leftarrow 0$
05. *do*
06. $D_i \leftarrow GP(f_{1..i})$
07. $fit_i \leftarrow SVM(D_i)$
08. $fitness \leftarrow fit_i$
09. $i = i + 1$
10. *While* ($fitness < Expected$) or ($i < |f|$)
11. $result \leftarrow \underset{j}{\operatorname{argmax}}(fit_j)$
12. $model \leftarrow SVM(D_{result})$

در ادامه، هر یک از بخش‌های گفته‌شده در الگوریتم ۱ به تفصیل تشریح می‌شود.

۳-۱- بالا بردن کیفیت داده

در الگوریتم‌های یادگیری ماشین فرض بر این است که داده‌ها فاقد نویز و همگی معتبر می‌باشند و هیچ نقصی در داده‌ها متصور نیست. این در حالی است که در دنیای واقعی داده‌ها خالی از نقص نیستند؛ بنابراین می‌توان بالا بردن کیفیت داده‌ها را جزئی جدایی‌ناپذیر از یک فرایند یادگیری ماشین موفق برشمرد. به همین منظور همان‌طور که در خط ۱ از الگوریتم ۱ قابل ملاحظه است، ابتدا مقادیر از دست‌رفته با مقادیری جایگزین شده است که امید می‌رود به مقدار واقعی و صحیح نزدیک باشد. برای این منظور داده‌های نامعلوم در هر ویژگی با شاخص پراکندگی مد جایگزین می‌شوند.

۳-۲- کاهش ویژگی مبتنی بر الگوریتم شبیه‌سازی تبرید

الگوریتم پیشنهادی کاهش ویژگی که معادل خط ۲ از الگوریتم ۱ است و به اختصار SAFR نامیده شده است، گامی ضروری در روش پیشنهادی است. همان‌طور که گفته شد، علت این ضرورت را می‌توان این‌گونه اظهار داشت که در خطوط ۵ تا ۱۰ از الگوریتم ۱، برای تغییر فضای مسئله و ایجاد بهترین مجموعه داده، امکان ساخت ژن‌هایی بر اساس تمامی ترکیبات ویژگی‌ها بررسی می‌شود؛ بنابراین با توجه به زمانبر بودن GP، این فاز به شدت زمان‌بر است؛ بنابراین باید ویژگی‌های بی‌ارزش تا حد امکان از چرخه کار کنار گذاشته شوند. الگوریتم‌های بهینه‌سازی سراسری برای این منظور با توجه به اینکه تنها بخشی از فضای مسئله

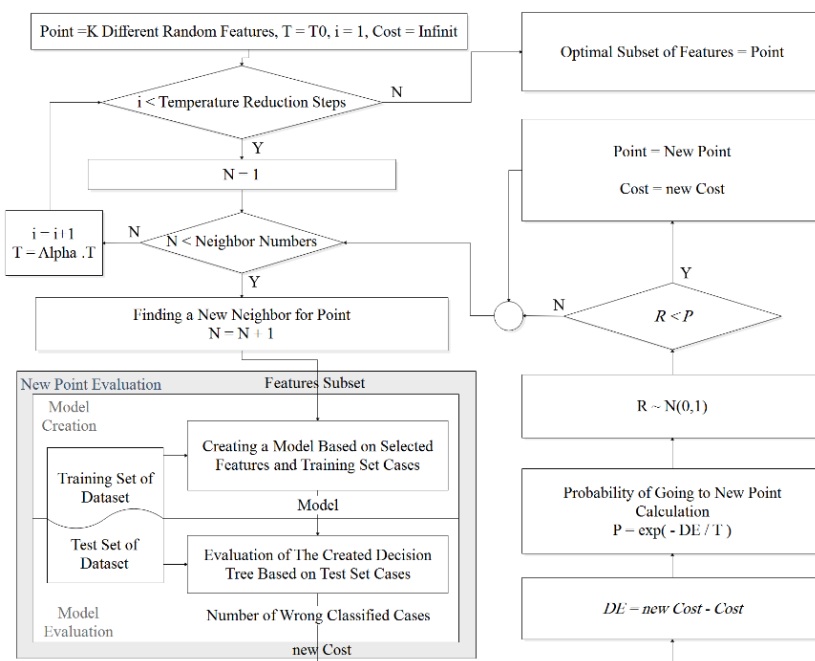
داده‌های آزمایشی داشته باشد. مدل‌سازی نقاط مربوط به الگوریتم شبیه‌سازی تبرید در مسئله به این صورت است که هر مجموعه k عضوی از ویژگی‌ها به‌عنوان یک نقطه از فضای جستجو در نظر گرفته می‌شود و

نقطه همسایه به نقطه‌ای اطلاق می‌شود که در یکی از ویژگی‌ها با نقطه جاری متفاوت باشد. نمونه‌ای از حالت $k=3$ در شکل ۲ آمده است.

الگوریتم ۲: کاهش ویژگی با الگوریتم شبیه‌سازی تبرید

```

01.  $T \leftarrow t_0$ 
02.  $\alpha \leftarrow \alpha_0$ 
03.  $Temp \leftarrow RandPerm(n, n)$ 
04.  $Train \leftarrow Temp_{1...2n/3}$ 
05.  $Test \leftarrow Temp_{2n/3+1...n}$ 
06. for  $i = 1$  to  $TempRedSteps$ 
07.      $FSS \leftarrow RandPerm(m, k)$ 
08.      $TrainData \leftarrow Data_{Train}.FSS$ 
09.      $TestData \leftarrow Data_{Test}.FSS$ 
10.      $model \leftarrow DecisionTree(TrainData)$ 
11.      $CurrentCost \leftarrow Eval(model, TestData)$ 
12.     for  $j = 1$  to  $Neighbors$ 
13.          $new \leftarrow newNeighbor(FSS)$ 
14.          $TrainData \leftarrow Data_{Train}.new$ 
15.          $TestData \leftarrow Data_{Test}.new$ 
16.          $model \leftarrow DecisionTree(TrainData)$ 
17.          $Cost \leftarrow Eval(model, TestData)$ 
18.          $\Delta E \leftarrow Cost - CurrentCost$ 
19.          $p \leftarrow e^{(-\Delta E/T)}$ 
20.          $n \leftarrow RandomN(0,1)$ 
21.         if  $n \leq P$ 
22.              $FSS \leftarrow new$ 
23.              $CurrentCost \leftarrow Cost$ 
24.     End
25.      $T \leftarrow \alpha.T$ 
26. End
27. return  $FSS$ 
    
```

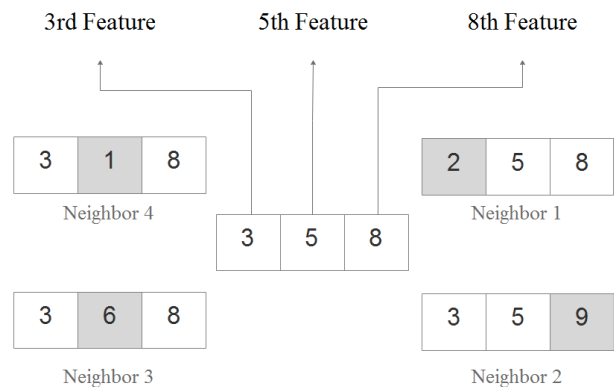


شکل ۱: روند کلی کاهش ویژگی با الگوریتم شبیه‌سازی تبرید

عنوان داده آزمایشی، مورد ارزیابی قرار می‌گیرد. در انتها درصد موارد طبقه‌بندی غلط در این ارزیابی به‌عنوان هزینه‌ی مربوط به پاسخ پیشنهادی الگوریتم شبیه‌سازی تبرید استفاده می‌شود. خطوط ۱۴ تا ۱۷ از الگوریتم ۲ روند محاسبه هزینه را نشان می‌دهد.

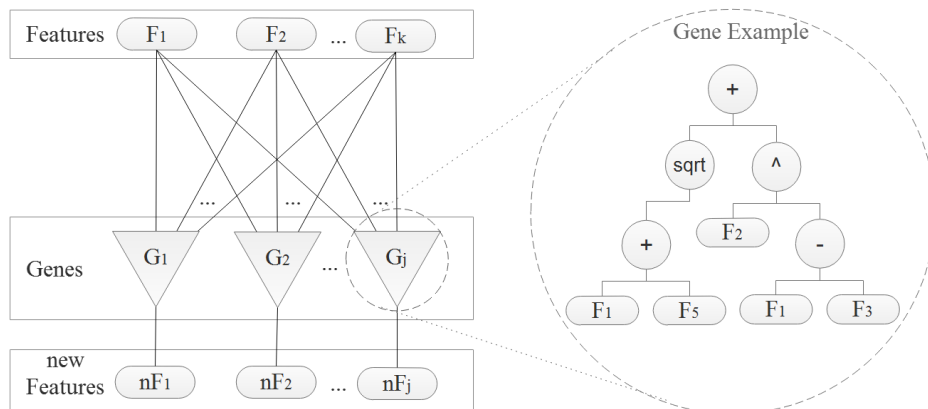
۳-۳- الگوریتم یادگیری ویژگی مبتنی بر برنامه‌نویسی ژنتیک

در بخش سوم که معادل خطوط ۵ تا ۱۰ از الگوریتم ۱ است، برای یادگیری ویژگی‌های جدید، یک الگوریتم برنامه‌نویسی ژنتیک ارائه شده است که فضای مسئله را از k به J تغییر می‌دهد. برای این منظور از تعداد J ژن برای یادگیری J ویژگی جدید استفاده می‌شود. روند کاری GP ارائه‌شده در شکل ۳ قابل مشاهده است. درواقع اگر مجموعه ویژگی‌های اولیه به‌صورت $\vec{F} = (F_1, F_2, \dots, F_k)$ در نظر گرفته شوند، عملکرد GP به‌صورت یک تابع GP، $n\vec{F} = GP(\vec{F})$ پیاده‌سازی شده است؛ که در آن nF بیانگر مجموعه ویژگی‌های جدید است. در انتها، در یک کار ابتکاری میزان برازش ژن‌ها با استفاده از الگوریتم SVM و بر اساس داده‌هایی انجام می‌شود که دربردارنده مقادیر ویژگی‌های جدید (nF) هستند.



شکل ۲: مدل‌سازی نقاط همسایگی در مسئله کاهش ویژگی با شبیه‌سازی تبرید

مدل‌سازی تابع هزینه نیز به این صورت است که در هر مرحله برای محاسبه هزینه‌ی مربوط به هر پاسخ الگوریتم شبیه‌سازی تبرید (که مشخص‌کننده k ویژگی است)، یک درخت تصمیم بر اساس زیرمجموعه‌ای از داده‌ها، شامل همان k ویژگی ایجاد می‌شود. مدل مذکور بر اساس درصدی از نمونه‌ها تحت عنوان داده آموزشی ایجاد می‌شود. پس از آن درخت ایجادشده بر روی بخشی از نمونه‌ها، تحت



شکل ۳: یادگیری ویژگی با برنامه‌نویسی ژنتیک

صحت: نسبت برچسب‌های درست به کل برچسب‌های زده‌شده

توسط مدل؛ یعنی:

$$Acc = \frac{TP+TN}{TP+TN+FP+FN} \quad (۶)$$

دقت: نسبت برچسب‌های مثبت درست تشخیص داده‌شده به کل

برچسب‌های مثبت تشخیص داده‌شده توسط مدل؛ یعنی:

$$Pr = \frac{TP}{TP+FP} \quad (۷)$$

بازیافت: نسبت تعداد موارد مثبت درست تشخیص داده‌شده به کل موارد

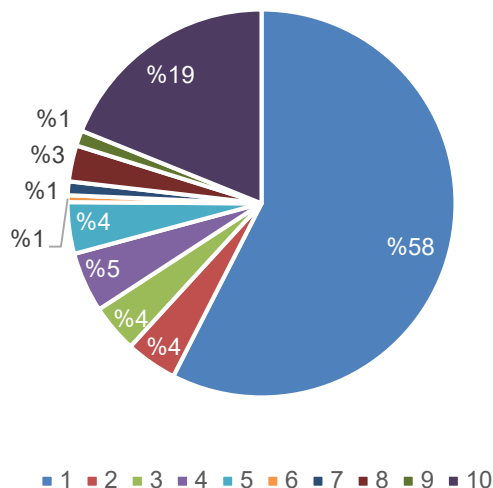
مثبت؛ یعنی:

$$Re = \frac{TP}{TP+FN} \quad (۸)$$

۴- ارزیابی

به‌منظور ارزیابی، روش پیشنهادی بر روی مجموعه داده‌های استاندارد، آزمایش شده و نتایج حاصل از آن با کارهای مشابه مقایسه شد. بدین منظور در زیربخش بعدی ابتدا مجموعه‌داده مورد استفاده تشریح شده و سپس نتایج حاصل از مقایسه ارائه می‌شود. در پیاده‌سازی روش پیشنهادی از یک سیستم با پردازنده دو هسته‌ای با توان پردازشی ۲/۱ گیگاهرتز استفاده شد که دارای ۱ گیگابایت حافظه اصلی است. پارامترهای اجرای الگوریتم شبیه‌سازی تبرید و GP نیز به ترتیب به‌صورتی تنظیم شده است که به ترتیب در جداول ۱ و ۲ دیده می‌شود. پارامترهای ارزیابی نیز شامل صحت، دقت، بازیافت و F_1 هستند که تعاریف و روش محاسبه هر کدام در ادامه بیان خواهد شد.

پراکندگی مد است. بر همین اساس، مقادیر ازدست‌رفته برای ویژگی مذکور با مقدار شاخص پراکندگی مد مربوط به همین ویژگی جایگزین شدند.



شکل ۴: نمودار توزیع مقادیر در ویژگی برهنگی هسته در مجموعه داده WBCD

۴-۲- نتایج و مقایسه

۴-۲-۱- نتایج تأثیر مراحل مختلف روش پیشنهادی

مقایسه تأثیر وجود و یا عدم وجود هر فاز بر دقت روش پیشنهادی می‌تواند ضرورت وجود فازهای مختلف را به اثبات برساند. این امر در جدول ۳ قابل مشاهده است.

همان‌طور که در جدول ۳ دیده می‌شود، کاهش ویژگی به تنهایی تعداد ویژگی‌ها را به میزان چشمگیری کاهش می‌دهد و نکته جالب توجه این است که کاهش ویژگی صورت گرفته همراه با افزایش دقت است. GP نیز به تنهایی ضمن کاهش چشمگیر ویژگی‌ها افزایش دقت را به دنبال دارد؛ اما این الگوریتم به شدت زمان‌بر است. ترکیب کاهش ویژگی و GP نیز کاهش بیش‌ازپیش ویژگی‌ها را ضمن افزایش دقت به دنبال دارد. البته این امکان متصور است که GP به تنهایی توان رسیدن به دقت روش ترکیبی پیشنهادی را با همان تعداد ویژگی‌ها داشته باشد؛ اما این امر مستلزم سپری شدن زمان بیش‌تری خواهد بود و این زمان بیش‌تر، در اجراهای مختلف متغیر بوده و می‌تواند در مواردی بسیار طولانی باشد. این درحالی است که استفاده از فاز کاهش ویژگی زمان مورد نیاز را برای رسیدن به پاسخ مطلوب، در حد تقریباً ثابتی حفظ خواهد کرد و این زمان ثابت، نسبت به زمان مورد نیاز برای ادامه GP عموماً کم‌تر است.

۴-۲-۲- تأثیر استفاده از الگوریتم‌های مختلف در فاز نهایی

استفاده از الگوریتم‌های مختلف در فاز اصلی یادگیری ماشین نیز، که فاز نهایی روش پیشنهادی است، در جدول ۴ آمده است. در این جدول نتیجه

معیار F1: ترکیبی است از معیارهای دقت و بازیافت که طبق رابطه x قابل محاسبه است. این معیار در واقع میانگین هارمونیک پارامترهای دقت و بازیابی است؛ یعنی:

$$F1 = 2 \times \frac{Pr \times Re}{Pr + Re} \quad (9)$$

جدول ۱: پارامترهای الگوریتم بهینه‌سازی شبیه‌سازی تبرید

پارامتر	مقدار
بیشترین تعداد تکرار	۲۰
بیشترین تعداد همسایگی	۵
نرخ کاهش دما	۰/۰۱
دمای اولیه	۱

جدول ۲: پارامترهای اولیه برنامه‌نویسی ژنتیک

پارامتر	مقدار
اندازه جمعیت	۱۰۰
تعداد نسل‌ها	۱۰۰
اندازه رقابت	۳
بیشینه عمق درخت	۵
بیشینه گره‌های درخت	بی‌نهایت
بیشینه تعداد ژن‌ها	۳
بازه اعداد ثابت	[۱۰ و -۱۰]

جدول ۳: مقایسه تأثیر وجود فازهای مختلف در روش پیشنهادی

	در مجموعه‌داده WBCD		
	Accuracy (%)	Feature space dimension	Time (s)
SVM	۹۷/۴۷	۳۰	۵/۷۶
SAFR+SVM	۹۷/۷۷	۱۳	۳۲۰۵/۵۶
GP+SVM	۹۷/۸۰	۸	۱۹۳۸۵/۶
SAFR+GP+SVM	۹۷/۸۹	۶	۲۲۴۹۸/۳۲

۴-۱- مجموعه‌داده

داده‌های به‌کاررفته در این مقاله شامل مجموعه‌داده استاندارد WBCD حاوی ۹ ویژگی از ۶۹۹ نمونه و مجموعه‌داده WDBC حاوی ۳۰ ویژگی از ۵۶۹ نمونه است. ارزیابی نتایج نیز بر روی هر دو مجموعه‌داده به صورت 10-Fold Cross Validation انجام شده است.

در مجموعه‌داده WBCD، ۱۶ مورد مقادیر ازدست‌رفته (Missing Value) وجود دارد و تمامی این ۱۶ مورد مربوط به ویژگی برهنگی هسته (bare nuclei) است؛ لذا حذف نمونه‌های دارای مقادیر ازدست‌رفته نمی‌تواند راه‌حل مناسبی باشد. همان‌طور که در شکل ۴ مشاهده می‌شود، نمودار توزیع مقادیر این ویژگی، مقادیری گسسته بین ۱ تا ۱۰ را نشان می‌دهد؛ بنابراین استفاده از شاخص پراکندگی مرکزی میانگین برای جایگزینی این مقادیر، انتخاب مناسبی نیست. با توجه به اینکه بیش از نیمی از مقادیر مربوط به ویژگی مذکور دارای مقدار عددی یک هستند واضح است که بهترین انتخاب برای مقادیر ازدست‌رفته شاخص

صورت گرفته است. در [۲۸] عمل طبقه‌بندی با استفاده از ۶ ویژگی انجام شده است. روش ابداعی در [۲۸] به این شکل بوده است که ابتدا به کمک الگوریتم K-means یک سری الگوی پنهان در بین داده‌ها یافته شده و به عنوان یک سری ویژگی جدید به مجموعه داده اضافه شده است. پس از آن مجموعه داده حاصل، به عنوان داده‌ی یادگیری به الگوریتم SVM داده شده و نتایج نیز به روش 10-Fold Cross Validation مورد ارزیابی قرار گرفته است. نتایج گویای این حقیقت است که این روش ابداعی توانسته دقتی ۹۷/۳۸٪ را به دنبال داشته باشد. مقایسه نتایج حاصل از روش پیشنهادی مقاله پیش رو بر روی مجموعه داده WDBC با آنچه که بر روی همین مجموعه داده در [۲۸] و [۲۹] انجام شده است در جدول ۶ قابل مشاهده است.

جدول ۶: مقایسه روش پیشنهادی با سایر روش‌های انجام شده بر

روى مجموعه داده WDBC		
	Feature space dimension	Accuracy (%)
ACO-SVM [29]	۱۵	۹۵/۹۶
GA-SVM [29]	۱۸	۹۷/۱۹
PSO-SVM [29]	۱۷	۹۷/۳۷
K-SVM [28]	۶	۹۷/۳۸
Proposed	۶	۹۷/۸۹

۴-۲-۴ - ارزیابی کلی نتایج روش پیشنهادی

نتایج ارزیابی روش پیشنهادی بر روی مجموعه داده WDBC در جدول ۷ قابل مشاهده است.

جدول ۷: ارزیابی آزمایش روش پیشنهادی روی مجموعه داده

WDBC				
	Acc	Re	Pr	F1
Benign	۹۵/۸۵	۹۸/۷۰	۹۷/۴۴	۹۸/۰۶
Malignant	۹۵/۸۵	۹۶/۹۲	۹۸/۴۳	۹۷/۶۷

نتایج حاصل از اعمال روش پیشنهادی بر روی مجموعه داده WBCD شامل ۹ ویژگی که مرتبط با بیماری سرطان سینه است نیز در جدول ۸ قابل ملاحظه است.

جدول ۸: ارزیابی آزمایش روش پیشنهادی روی مجموعه داده

WBCD				
	Acc	Re	Pr	F1
Benign	۹۵/۸۵	۹۶/۵۱	۹۷/۱۴	۹۶/۸۲
Malignant	۹۵/۸۵	۹۴/۶۱	۹۳/۴۴	۹۴/۰۲

۵- بحث

فلسفه روش پیشنهادی بر دو پایه استوار است. اول اینکه بالا بودن تعداد ویژگی‌ها نه تنها لزوماً منجر به دقت بالا نمی‌شود؛ بلکه در برخی موارد منجر به افت دقت نیز خواهد شد؛ بنابراین کاهش ویژگی در چنین

حاصل از اجرای الگوریتم‌های مرسوم یادگیری ماشین به عنوان فاز نهایی در روش پیشنهادی به نمایش درآمده است. این امر استفاده از الگوریتم SVM را در روش پیشنهادی به وضوح توجیه می‌کند.

جدول ۴: مقایسه الگوریتم‌های مرسوم یادگیری به عنوان فاز

نهایی روش پیشنهادی در مجموعه داده WDBC		
	Accuracy (%)	Time (ms)
Decision Tree	۹۶/۹۶	۳۷۳
KNN	۹۶/۷۰	۱۹۰
Neural Networks	۹۵/۱۴	۱۱۰۲۱
SVM	۹۷/۸۹	۵۶۰

۴-۲-۳ - مقایسه نتایج با کارهای مشابه

مقایسه این نتایج با کارهایی که بر روی مجموعه داده WBCD در [۲۳] و [۲۴] صورت پذیرفته است افزایش دقت را ضمن کاهش چشمگیر ویژگی‌ها نشان می‌دهد. به گونه‌ای که تنها با استفاده از دو ویژگی، عملکرد مذکور را به دست آورده است. مقایسه نتایج حاصل از روش پیشنهادی با کارهای انجام شده بر روی مجموعه داده WBCD در جدول ۵ قابل مشاهده است.

جدول ۵: مقایسه روش پیشنهادی با سایر روش‌های انجام شده در

مجموعه داده WBCD		
Methods		Accuracy (%)
Karabatak [23]	AR+ ANN	۹۵/۶۰
Subashini [24]	SVM+RBFNN	۹۵/۵۶
Proposed	SAFR+GP+ SVM	۹۵/۸۵

در [۲۳] که بر روی مجموعه داده WBCD صورت گرفته است. ابتدا به کمک AR ابعاد مجموعه داده کاهش یافته است؛ به صورتی که تعداد ویژگی‌ها از ۹ ویژگی به ۴ ویژگی رسیده است. در ادامه، مجموعه داده حاصل از کاهش ویژگی به الگوریتم شبکه عصبی داده شد و نتایج به روش 3-Fold Cross Validation ارزیابی شده است. نتایج نشان می‌دهد روش پیشنهادی آن‌ها، عمل طبقه‌بندی را با دقت ۹۵/۶٪ انجام داده است. در [۲۴] نیز که بر روی همین مجموعه داده، شامل ۹ ویژگی، صورت گرفته است؛ نتایج بیانگر این حقیقت است که RBFNN و SVM در مسئله مذکور به ترتیب دقتی برابر با ۹۵/۵۶٪ و ۹۲/۱۳٪ داشته‌اند.

روش ارائه شده در [۲۸] به منظور افزایش دقت طبقه‌بندی، بر روی مجموعه داده WDBC مورد ارزیابی قرار گرفته است. بر همین اساس، روش پیشنهادی مقاله پیش رو با [۲۸] مقایسه شده است. اعمال روش پیشنهادی بر روی داده‌های WDBC، شامل ۳۰ ویژگی مربوط به سرطان سینه، و ارزیابی نتایج به صورت 10-Fold Cross Validation منجر به تشخیص بدخیمی و خوش‌خیمی با دقت ۹۷/۸۹٪ شده است که در مقایسه با روش پیشنهادی در [۲۸] دقت بالاتری است و این افزایش دقت، ضمن کاهش ویژگی‌ها دقیقاً به اندازه کاری است که در [۲۸]

مراجع

- [1] "Institute for Research, Education and Treatment of Cancer", www.ncii.ir, [Online Access: Feb. 22, 2015].
- [2] "Breastcancer.org", www.breastcancer.org, [Online Access: Feb. 22, 2015].
- [3] S. Kim, Y. Choi, M. Lee, "Deep learning with support vector data description," *Neurocomputing*, vol. 165, no. 1, pp. 111-117, Oct. 2015.
- [4] A.H. Gandomi, D. Mohammadzadeh, J.L. Pérez-Ordóñez, A.H. Alavi, "Linear genetic programming for shear strength prediction of reinforced concrete beams without stirrups," *Applied Soft Computing*, vol. 19, no. 1, pp. 112-120, Feb. 2014.
- [5] Sh. Ding, L. Lin, G. Wang, H. Chao, "Deep Feature Learning with Relative Distance Comparison for Person Re-identification," *Pattern Recognition*, vol. 48, no. 1, pp. 2993-3003, Oct. 2015.
- [6] E. Emary, H.M. Zawbaa, A.E. Hassanien, "Binary grey wolf optimization approaches for feature selection," *Neurocomputing*, vol. 172, no. 1, pp. 371-381, Jan. 2016.
- [7] I. Zyout, J. Czajkowska, M. Grzegorzec, "Multi-scale textural feature extraction and particle swarm optimization based model selection for false positive reduction in mammography," *Computerized Medical Imaging and Graphics*, vol. 46, pp. 95-107, Feb. 2015.
- [8] H. Min, R. Weijie, "Global mutual information-based feature selection approach using single-objective and multi objective optimization," *Neurocomputing*, vol. 168, pp. 47-54, Nov. 2015.
- [9] K. Shima, E. Mahdi, "Feature selection using multimodal optimization techniques," *Neurocomputing*, vol. 171, pp. 586-597, Jan. 2015.
- [10] M. Parham, R. Mehrdad, "Integration of graph clustering with ant colony optimization for feature selection," *Knowledge-Based Systems*, vol. 84, pp. 144-161, Aug. 2015.
- [11] N.K. Sreeja, A. Sankar, "Pattern Matching based Classification using Ant Colony Optimization based Feature Selection," *Applied Soft Computing*, vol. 31, pp. 91-102, Jun. 2015.
- [12] A.T. Mohammad Reza, H. Seyyed Abed, N.S. Mohammad Bagher, "Evaluation of Visual Selective Attention by Event Related Potential Analysis in Brain Activity," *Tabriz Journal of Electrical Eng*, vol. 46, no. 1, pp. 13-24, 2016.
- [13] B. Morteza, P. Hosein, "Epilepsy Detection Based on Optimization of Fused Hartley Transform Feature with Hybrid Model of MLP and GA using Memetic Learning Strategy," *Tabriz Journal of Electrical Eng*, vol. 45, no. 4, pp. 51-67, 2015.
- [14] J. Yu, J. Yu, A.A. Almal, S.M. Dhanasekaran, D.Ghosh, W.P. Worzelz, A.M. Chinnaiyan, "Feature Selection and Molecular Classification of Cancer Using Genetic Programming," *Neoplasia*, vol. 9, no. 4, pp. 292-303, April. 2007
- [15] F. Archetti, I. Giordani, L. Vanneschi, "Genetic programming for anticancer therapeutic response prediction using the NCI-60 dataset," *Computers & Operations Research*, vol. 37, no. 1, pp. 1395-1405, Mar. 2009.
- [16] G. Phillips-Wren, P. Sharkey, S. Morss Dy, "Mining lung cancer patient data to assess healthcare resource utilization," *Expert Systems with Applications*, vol. 35, no. 1, pp. 1611-1619, Aug. 2008.
- [17] L. Li, H. Tang, Z. Wu, J. Gong, M. Gruidl, J. Zou, M. Tockman, R.A. Clark, "Data mining techniques for cancer detection using serum proteomic profiling," *Artificial Intelligence In Medicin*, vol. 32, no. 1, pp. 71-83, Mar. 2004.
- [18] S.C. Shah, A. Kusiak, M.A. O'Donnell, "Patient-recognition data-mining model for BCG plus interferon immunotherapy bladder cancer treatment," *Computers in Biology and Medicine*, vol. 36, no. 1, pp. 634-655, Mar. 2005.
- [19] K.R. Seeja, "Feature selection based on closed frequent itemset mining: A case study on SAGE data classification," *Neurocomputing*, vol. 151, pp. 1027-1032, Mar. 2015.
- [20] S. Shah, . Kusiak, "Cancer gene search with data-mining and genetic algorithms," *Computers in Biology and Medicine*, vol. 37, no. 1, pp. 251-261, Jan. 2007.
- [21] Y. Qiang, Y. Guo, X. Li, Q. Wang, H. Chen, D. Cuic, "The Diagnostic Rules of Peripheral Lung Cancer Preliminary Study

مواردی می‌تواند دقت را نیز افزایش دهد. دلیل این امر، وجود ویژگی‌های بی‌فایده در کنار ویژگی‌های مخرب است [۳۰]؛ بنابراین کاهش ویژگی می‌تواند با حذف چنین ویژگی‌هایی افزایش دقت طبقه‌بندی را به دنبال داشته باشد. دوم این که عملیات یادگیری ویژگی بر روی داده‌های اولیه و خام، می‌تواند منجر به ایجاد ویژگی‌های سطح بالا و متمایزکننده جدید شود؛ به‌گونه‌ای که ویژگی‌های خروجی الگوریتم یادگیری ویژگی، دقت بالاتری برای طبقه‌بندی داشته باشد. استفاده از الگوریتم GP نشان داد که این الگوریتم توانایی بالایی در یادگیری ویژگی دارد. گرچه در بطن این الگوریتم، عمل کاهش ویژگی نیز نهفته است؛ اما کاهش ویژگی اولیه (که در این پژوهش با الگوریتم بهینه‌سازی شبیه‌سازی تبرید انجام شده است) می‌تواند منجر به تسریع و بهبود الگوریتم GP شود.

نکته قابل ذکر، کاربرد روش ارائه‌شده در امور بالینی است. مجموعه‌داده‌های استفاده‌شده در این پژوهش، توسط مراکز بالینی جمع‌آوری شده است و طبقه‌بندی خوش‌خیمی و بدخیمی نیز توسط متخصصین مربوطه انجام شده است. از آنجایی که روش ارائه‌شده یک روش یادگیری تحت نظارت است، نتایج به‌دست‌آمده نشان می‌دهد که مدل حاصل از روش ارائه‌شده، توانسته با دقت بالایی به عملکرد متخصصین نزدیک شود.

نکته دیگری که در این پژوهش باید به درستی مورد بحث قرار گیرد، بالا بودن زمان یادگیری است. باید توجه داشت که این زمان، صرف یادگیری در فاز آفلاین خواهد شد؛ لذا در صورت استفاده عملی، تفاوتی در سرعت سیستم تشخیصی مبتنی بر روش پیشنهادی، نسبت به سایر روش‌ها وجود نخواهد داشت.

۶- نتیجه‌گیری

بر کسی پوشیده نیست که GP الگوریتمی مناسب برای برازش منحنی و رگرسیون است اما همان‌طور که در این مقاله اثبات شد، الگوریتم مذکور قابلیت استفاده در مسائل یادگیری ویژگی را نیز به‌خوبی دارد. اگرچه در این مقاله ابتدا تعداد ویژگی‌ها به‌کمک الگوریتم شبیه‌سازی تبرید کاهش یافت و این امر بیانگر قدرت الگوریتم‌های فراابتکاری و به‌طور خاص الگوریتم شبیه‌سازی تبرید در کاهش ویژگی است؛ اما نکته جالب‌توجه در این پژوهش این است که ژن‌های ایجادشده توسط الگوریتم GP تنها شامل تعدادی از ویژگی‌های کاهش‌یافته است و این حکایت از این حقیقت دارد که الگوریتم GP، عمل کاهش ویژگی را نیز به‌صورت ضمنی انجام می‌دهد و این به معنی کاهش بیش‌ازپیش ابعاد مسئله است. البته در صورت حذف فاز کاهش ویژگی، رسیدن به نتیجه مطلوب با استفاده از الگوریتم GP زمان بیش‌تری را می‌طلبد. با توجه به زمان‌بر بودن GP، فاز کاهش ویژگی جزئی لازم و جدایی‌ناپذیر در روش پیشنهادی است. نتایج این مقاله بیانگر این حقیقت است که الگوریتم GP می‌تواند به‌منظور یادگیری ویژگی در مسائلی به کار برده شود که تعداد ویژگی‌های اولیه بسیار زیادی دارند.

- based segmentation and multilayer perceptron neural networks classifier," *Applied Soft Computing*, vol. 14, no. 1, pp. 62-71, Aug. 2013.
- [27] P.J. García-Laencina, P.H. Abreu, M.H. Abreu, N. Afonso, "Missing data imputation on the 5-year survival prediction of breast cancer patients with unknown discrete values," *Computers in Biology and Medicine*, vol. 59, no. 1, pp. 125-133, Feb. 2015.
- [28] B. Zheng, S.W. Yoon, S.S. Lam, "Breast cancer diagnosis based on feature extraction using a hybrid of K-means and support vector machine algorithms," *Expert Systems with Applications*, vol. 41, no. 1, pp. 1476-1482, March. 2014.
- [29] Prasad. Y, Biswas. k, Jain. c, "Svm classifier based feature selection using ga, aco and pso for sirna design," *In Proceedings of the first international conference on advances in swarm intelligence*, pp. 307-314, Beijing, China, June. 2010.
- [30] MRC Human Genetics Unit, MRC Institute of Genetics and Molecular Medicine, University of Edinburgh et al, " Application of high-dimensional feature selection: evaluation for genomic prediction in man," *Scientific Reports*, vol. 5, pp. 1-12, May. 2015.
- Based on Data Mining Technique," *Journal of Nanjing Medical University*, vol. 21, no. 3, pp. 190195, Nov. 2006.
- [22] T.Z. Tan, C. Quek, G.S. Ng, E.Y.K. Ng, "A novel cognitive interpretation of breast cancer thermography with complementary learning fuzzy neural memory structure," *Expert Systems with Applications*, vol. 33, no. 1, pp. 652-666, Jun. 2007.
- [23] M. Karabatak, M.C. Ince, "An expert system for detection of breast cancer based on association rules and neural network," *Expert Systems with Applications*, vol. 36, no. 1, pp. 3465-3469, Feb. 2009.
- [24] T.S. Subashini, V. Ramalingam, S. Palanivel, "Breast mass classification based on cytological patterns using RBFNN and SVM," *Expert Systems with Applications*, vol. 36, no. 1, pp. 5284-5290, Jun. 2009.
- [25] M. Kowa, P. Filipczuk, A. Obuchowicz, J. Korbicz, R. Monczak, "Computer-aided diagnosis of breast cancer based on fine needle biopsy microscopic images," *Computers in Biology and Medicine*, vol. 43, no. 1, pp. 1563-1572, Aug. 2013.
- [26] A.E.Hassaniena, H.M. Moftah, A.T. Azar, M. Shoman, "MRI breast cancer diagnosis hybrid approach using adaptive ant-

زیر نویس ها

\ Simulated Annealing Feature Reduction