

Diagnosis of Breast Cancer Subtypes using the Selection of Effective Genes from Microarray Data

Tabatabaei A¹, Derhami V^{1*}, Sheikhpour R², Pajoohan M-R¹

¹ Computer Engineering Department, Faculty of Engineering, Yazd University, Yazd, Iran

² Department of Computer Engineering, Faculty of Engineering, Ardakan University, Ardakan, Iran

Receive: 2019/2/13

Accepted: 2019/4/14

*Corresponding Author:

Vali Derhami

vderhami@yazd.ac.ir

Ethics Approval:

Abstract

Introduction: Early diagnosis of breast cancer and the identification of effective genes are important issues in the treatment and survival of the patients. Gene expression data obtained using DNA microarray in combination with machine learning algorithms can provide new and intelligent methods for diagnosis of breast cancer.

Methods: Data on the expression of 9216 genes from 84 patients across 5 different types of cancer was obtained using microarray technology. In this study, we proposed a feature selection method based on the correlation between abnormal expression of genes and cancer for diagnosis of breast cancer. Then, we used K-nearest neighbor (KNN), support vector machine (SVM), and naive Bayesian (NB) classifiers to evaluate the performance of the proposed method in the selection of relevant genes.

Results: The proposed feature selection method coupled with the KNN classifier predicted all types of cancer with 100% accuracy and using 38 of the 9216 genes. The proposed method could also identify the genes associated with each class. Moreover, the proposed feature selection method coupled with NB and SVM classifiers achieved accuracy rates of 90% and 96.67% using 17 and 22 genes, respectively.

Conclusion: The results of this study demonstrated that the proposed feature selection method has better performance compared with other methods. The proposed method is able to distinguish the genes involved in each cancer class and detect overexpression or underexpression of selected genes, which can be used by physicians and researchers in the field of health care.

Keywords: Breast Cancer, Feature Selection, Microarray Data, Classification

تشخیص نوع سرطان پستان با استفاده از انتخاب ژن‌های موثر از داده‌های ریزآرایه

سید ابوالفضل طباطبایی^۱، ولی درهمی^{۱*}، راضیه شیخ‌پور^۲، محمد رضا پژوهان^۱

^۱ گروه مهندسی کامپیوتر، پردیس فنی و مهندسی، دانشگاه یزد، یزد، ایران

^۲ گروه مهندسی کامپیوتر، دانشکده فنی و مهندسی، دانشگاه اردکان، اردکان، ایران

چکیده

تاریخ ارسال: ۹۷/۱۱/۲۴

تاریخ پذیرش: ۹۸/۱/۲۵

نشانی نویسنده مسئول:

ولی درهمی

vderhami@yazd.ac.ir

مقدمه: تشخیص زودهنگام سرطان پستان و ژن‌های موثر در آن نقش بسیار کلیدی در درمان و حیات بیمار ایفا می‌کند. با استفاده از داده‌های بیان ژن استخراج شده از فناوری ریزآرایه و الگوریتم‌های یادگیری ماشین می‌توان روش‌های نوین و هوشمندی در نظام سلامت و درمان ارایه داد که با دقت بالایی قادر به تشخیص سرطان پستان باشند.

روش بررسی: داده‌های استفاده شده در این پژوهش، شامل داده‌های بیان ۹۲۱۶ ژن مربوط به ۸۴ بیمار در ۵ نوع مختلف سرطان است که با استفاده از فناوری ریزآرایه به دست آمده است. در این مطالعه برای افزایش کارایی سیستم‌های تشخیص سرطان پستان، روش انتخاب ویژگی مبتنی بر ارتباط بین بیان غیرنرم‌مال ژن و سرطان ارایه شده است. سپس از سه دسته‌بند پرکاربرد و رایج K-نزدیکترین همسایه (KNN)، ماشین بردار پشتیبان (SVM) و بیزی ساده (NB) برای سنجش کارایی ژن‌های انتخاب شده استفاده شد.

یافته‌ها: بررسی‌های انجام شده نشان دادند که با روش انتخاب ویژگی پیشنهادی و با استفاده از دسته‌بند KNN می‌توان فقط با انتخاب ۳۸ ژن از میان ۹۲۱۶ ژن مربوط به داده‌های بیماران سرطانی که با استفاده از فناوری ریزآرایه به دست آمده است، انواع سرطان پستان بیماران در داده‌های آزمایش را با دقت ۱۰۰٪ تشخیص داد و ژن‌های مرتبط با هر کلاس را نیز تفکیک کرد. همچنین با استفاده از دسته‌بند NB، دقت ۹۰٪ با ۱۷ ژن و با دسته‌بند SVM، دقت ۹۶/۶۷٪ با ۲۲ ژن انتخابی به دست آمد.

نتیجه‌گیری: نتایج این مطالعه نشان دادند که روش انتخاب ویژگی پیشنهادی، با در نظر گرفتن همزمان دو معیار دقت و تعداد ژن انتخابی، عملکرد مناسبی نسبت به سایر روش‌ها دارد. توانایی تفکیک ژن‌های موثر در هر کلاس سرطان، به علاوه تشخیص بیان بیشتر از حد یا کمتر از حد ژن‌های انتخابی، از خصوصیات ویژه روش پیشنهادی است که می‌تواند مورد استفاده متخصصین و پژوهشگران حوزه درمان و مراقبت قرار گیرد.

واژه‌های کلیدی: سرطان پستان، انتخاب ویژگی، داده‌های ریزآرایه، دسته‌بندی

بیماری است. مطالعات و روش‌های سنتی، تعداد زیادی از زن‌ها را به عنوان زن‌های مؤثر در بیماری در نظر می‌گیرند (۱۳). به دلیل این‌که روش‌های تجربی برای تشخیص زن‌های مؤثر از میان تعداد زیاد زن‌های نامزد بسیار هزینه‌براست، به کارگیری روش‌های انتخاب ویژگی (زن) توصیه می‌شود. به علاوه وجود نویز^۹ و تعداد بسیار بالای اطلاعات زنی، تحلیل داده‌های ریزآرایه را به یک دامنه مهیج تبدل کرده است (۱۴). مطالعات متعددی نشان داده‌اند که تعداد بسیاری از داده‌های بیان زن ریزآرایه DNA (ویژگی‌ها) از نظر معیارهای دسته‌بندی دارای بار اطلاعاتی نیستند (داده‌ی غیرمرتبط) (۱۵). بنابراین انتخاب ویژگی (زن) و فرآیند تشخیص و حذف ویژگی‌های غیرمرتبط، نقشی حیاتی در تحلیل داده‌های ریزآرایه استخراج شده از DNA ایفا می‌کنند.

یافته‌های علمی به خصوص طی دو دهه اخیر ارتباط بین بیان غیرنرم‌مال زن و انواع سرطان‌ها را نشان می‌دهد (۱۶-۱۸). جهش^{۱۰}، بیان بیش از حد^{۱۱} و بیان کمتر از حد^{۱۲} زن مصادیق بیان غیرنرم‌مال زن هستند. مارتین کورنا^{۱۳} و همکاران در (۱۹) ارتباط بین جهش سلوی با برخی سرطان‌ها را نشان دادند. البته باید توجه داشت که هر جهش منجر به سرطان نمی‌شود. زین العابدین و همکاران در (۲۰) دسته‌ای از زن‌های با بیان کمتر از حد و دسته‌ای از زن‌ها با بیان بیش از حد را برای تشخیص نوعی سرطان پوست به نام ملانوما^{۱۴} معرفی کردند. به نظر می‌رسد می‌توان از ارتباط بین بیان زن و سرطان به عنوان دانش و به صورت مهندسی معکوس در مسئله انتخاب زن در داده‌های ریزآرایه استفاده نمود که منجر به کاهش پیچیدگی مسئله خواهد شد.

مواد و روش‌ها

مطالعه حاضر توصیفی و داده محور است که به ارایه روشی برای تشخیص انواع سرطان پستان از دید مولکولی با استفاده از داده‌های بیان زن بیماران سرطانی پرداخته است.

مقدمه

سرطان پستان شایع‌ترین سرطان و دومین عامل مرگ و میر در زنان آمریکایی است (۱). در ایران سالیانه بیش از ۷۰۰۰ مورد به مبتلایان سرطان پستان اضافه می‌شود (۲). در این بیماری سلول‌های بدخیم (دارای رشد غیر قابل کنترل) در بافت پستان شروع به فعالیت می‌کنند (۳، ۴). علاوه بر پارامترهای تشخیص سنتی بیماری مثل اندازه و درجه (گرید) تومور، معمولاً از وضعیت مثبت یا منفی بودن نشان‌گرهای ایمونوھیستوشیمی^۱ IHC نظیر گیرنده‌های^۲ ER و^۳ PR و^۴ HER2 و نشان‌گر^۵ Ki67 برای تشخیص و مراقبت از بیماران سرطانی استفاده می‌شود (۵، ۶). از دیدگاه مولکولی، انواع تومورهای پستان به پنج دسته اصلی با علایم بالینی متمایز تقسیم می‌شوند (۷، ۸). این دسته‌ها عبارتند از: Luminal A, Luminal B, HER2 over-expression, Basal-Normal-like, Normal-like. غیر از گروه سایر گروه‌ها از نشان‌گرهای ایمونوھیستوشیمی متمایزی برخوردار هستند (۹). به عنوان مثال در گروه Basal-like، نشان‌گرهای ER و PR و HER2 هر سه منفی هستند.

یکی از راه‌های تشخیص انواع سرطان، استفاده از داده‌های بیان زن بیماران است. پیشرفت فناوری در حوزه بیوانفورماتیک، به خصوص فناوری ریزآرایه^۶ باعث شده که داده‌های بیان زنی^۷ هزاران زن مربوط به یک نمونه مبتلا به سرطان به طور همزمان استخراج شود (۱۰). البته این استخراج داده مستلزم هزینه‌های بالا است، به همین دلیل تعداد نمونه‌ها در مقایسه با تعداد ویژگی‌های^۸ استخراج شده بسیار کم است. چنین داده‌هایی به اختصار HDLSS^۹ نامیده می‌شود (۱۱) که در آن معمولاً تعداد نمونه‌ها کمتر از ۱۰۰ و تعداد ویژگی‌ها بین ۶۰۰۰ تا ۶۰۰۰۰ است (۱۲). تشخیص زن‌های مؤثر در بیماری از داده‌های ریزآرایه، یک موضوع مهم در رابطه با ارتقا دانش در مورد مکانیزم بیماری و بهبود روش‌های مواجهه با

^۱ Immunohistochemistry

^۲ Estrogen receptor

^۳ Progesterone receptor

^۴ Human epidermal growth factor receptor 2

^۵ Micro Array Data set

^۶ Gene expression data

^۷ Feature

^۸ High Dimension Low Sample Size

^۹ Noise

^{۱۰} Mutation

^{۱۱} Overexpression

^{۱۲} Underexpression or misexpression or less expression

^{۱۳} Martincorena

^{۱۴} Melanoma

مقادیر ز از ۱ تا ۵ تغییر می‌کند. $\text{Score}(f_i|C_j)$ میزان اهمیت و امتیاز ویژگی (زن) \mathfrak{z} ام در کلاس \mathfrak{z} ام را نشان می‌دهد. α . تعداد افزار^{۱۶} کلاس C_j توسط ویژگی f_i است. منظور از افزار یک مجموعه، تقسیم آن به تعدادی زیرمجموعه است که اشتراک آن‌ها تهی و اجتماع آن‌ها برابر مجموعه اولیه باشد (۲۱). هریک از زیربخش‌های کلاس C_j شامل نمونه‌هایی است که از نظر مقداری کنار هم قرار دارند. L_k ، تعداد نمونه‌های زیر بخش k ام از زیربخش‌های افزارشده کلاس C_j است. هرچه عدد α کوچک‌تر باشد، کلاس C_j به زیربخش‌های بزرگ‌تر تقسیم می‌شود که به معنی وجود نمونه‌های بیشتر با خصوصیات مشترک است. بنابراین، طبق رابطه (۱) امتیاز بیشتری به ویژگی f_i تعلق می‌گیرد. $|C_j|$ نیز معرف تعداد نمونه‌های کلاس \mathfrak{z} ام است.

انتخاب زن‌ها براساس امتیاز به‌دست آمده: بعد از مشخص شدن امتیاز ویژگی‌ها در هر کلاس، نوبت تصمیم‌گیری در مورد ترتیب نهایی زن‌ها برای استفاده در مدل‌های رایج یادگیری مثل^{۱۷} SVM و^{۱۸} KNN و^{۱۹} NB است. یکی از ساده‌ترین راه حل‌ها این است که ویژگی‌های با ارزش بالاتر از هر کلاس به ترتیب در لیست خروجی قرار گیرند. در این صورت تاثیر ویژگی در سایر کلاس‌ها در نظر گرفته نمی‌شود. راه حل پیشنهادی در این مطالعه این است که با استفاده از رابطه (۲) امتیاز یک ویژگی در سایر کلاس‌ها نیز در محاسبه امتیاز نهایی ویژگی در کلاس مورد نظر لحاظ گردد و سپس ویژگی‌های با ارزش بالاتر از هر کلاس به ترتیب در لیست خروجی قرار گیرند.

$$\text{Score}(f_i|C_j) = \text{Score}(f_i|C_j) + \beta * \sum_{r \neq j} \text{Score}(f_i|C_r) \quad (2)$$

ضریب β میزان مشارکت سایر کلاس‌ها در امتیازدهی را مشخص می‌کند. مقدار β باید طوری انتخاب شود که تاثیر و نقش ویژگی در کلاس اصلی کم‌رنگ شود. از مقادیر منفی β ، برای به‌دست آوردن ویژگی‌های متمایز و از مقادیر مثبت آن، برای به‌دست آوردن ویژگی‌های مشترک بین کلاس‌ها استفاده می‌شود. انتخاب مقدار

تصویف مجموعه داده‌ها: داده‌های استفاده شده در این مطالعه، شامل داده‌های بیان ۹۲۱۶ زن مربوط به ۸۴ بیمار سلطانی است که توسط پراو^{۱۵} و همکاران با استفاده از فناوری ریزآرایه به دست آمده است (۷). داده‌ها شامل پنج کلاس (نوع) مختلف سلطان می‌باشد که کلاس ۱ با ERBB2 برچسب Luminal، کلاس ۲ با برچسب Basal-like کلاس ۳ با برچسب Cell_lines و کلاس ۵ با برچسب Normal-like مشخص شده است. مجموعه داده‌های مذکور به صورت دو دسته آموزش و آزمون در دسترس هستند که در ۵۴ نمونه‌ی دسته آموزش به ترتیب ۲۰، ۶، ۷ و ۱۲ نمونه از هر کلاس وجود دارد و در ۳۰ نمونه‌ی مربوط به دسته آزمون، نمونه‌ها به ترتیب ۱۲، ۳، ۵ و ۷ در کلاس‌های نامبرده قرار دارند.

روش پیشنهادی: در این بخش، از ارتباط بین بیان غیرنرم‌الzen و سلطان به عنوان دانش، برای شناسایی و انتخاب زن‌های موثر در تشخیص سلطان پستان استفاده می‌شود و یک روش انتخاب ویژگی برای انتخاب زن‌های موثر در داده‌های ریزآرایه سلطان پستان پیشنهاد می‌شود که شامل دو بخش امتیازدهی زن‌ها در هر کلاس و انتخاب زن‌ها براساس امتیاز به‌دست آمده در هر کلاس است.

امتیازدهی زن‌ها در هر کلاس: داده‌های بیان زن ریزآرایه معرف کمیت توده mRNA مرتبط با نمونه هر یک از بیماران است (۱۲)، بنابراین هر تغییر در بیان زن عامل بیماری باعث می‌شود مقادیر آن زن نیز از حالت نرم‌الzen خارج شود. در نتیجه، مقادیر آن زن در نمونه‌های مختلف از یک بیماری به‌هم نزدیک می‌شوند و در کنار هم یک پیوستگی محدود را در کلاس مربوط به آن بیماری تشکیل دهند. در این مطالعه برای کمیت‌سازی این پیوستگی‌ها رابطه زیر پیشنهاد می‌شود:

$$\text{Score}(f_i|C_j) = \sum_{k=1}^{\alpha} \frac{L_k}{|C_j|} \exp\left(\frac{L_k}{|C_j|}\right) \quad (1)$$

که f_i ویژگی مربوط به بیان زن \mathfrak{z} ام است. با توجه به تعداد زن‌ها، مقادیر i از ۱ تا ۹۲۱۶ متغیر است. C_j ، کلاس \mathfrak{z} ام سلطان را نشان می‌دهد که با توجه به ۵ کلاس سلطان،

¹⁶ Partition

¹⁷ Support Vector Machine

¹⁸ K-nearest neighbor

¹⁹ Naive Bayes

¹⁵ Peru

جدول ۱: عملکرد روش پیشنهادی با توجه به معیار دقت

$\beta = -0.2$	$\beta = 1$	$\beta = 0.01$	$\beta = 0$	دسته‌بند
۸۳/۳۳	۹۰	۹۰	۸۶/۶۷	NB
۹۳/۳۳	۸۶/۶۷	۹۶/۶۷	۹۶/۶۷	SVM
۸۳/۳۳	۸۶/۶۷	۱۰۰	۹۶/۶۷	KNN

جدول ۲: عملکرد روش پیشنهادی با توجه به معیار تعداد ژن

$\beta = -0.2$	$\beta = 1$	$\beta = 0.01$	$\beta = 0$	دسته‌بند
۲۶	۱۷	۳۹	۶۰	NB
۴۱	۳۰	۲۲	۴۵	SVM
۳۰	۵۱	۳۸	۴۳	KNN

توجه جدول‌های (۱) و (۲)، در مجموع مدل‌های یادگیری SVM و KNN روی ژن‌های انتخابی عملکرد بهتری نسبت به مدل یادگیری NB دارند. در بهترین وضعیت با استفاده از مدل یادگیری KNN و با تعداد ۳۸ ژن انتخابی، نوع سرطان در کلیه نمونه‌های آزمایش به درستی تشخیص داده شد. همچنین دقت قابل قبول ۹۶/۶۷ با تعداد ۲۲ ژن در مدل SVM بدست آمد. جدول‌های (۳) و (۴) کارایی روش پیشنهادی را در مقایسه با سه روش انتخاب ویژگی معروف Relief (۲۲) و ARCO (۲۳) و FAST (۲۴) در شرایط یکسان (ماکریم ۱۰۰ ژن انتخابی و دسته‌بندهای مشابه) نشان می‌دهد. مقایسه نتایج بر اساس میانگین خطاهای در هر کلاس صورت گرفته است.

با توجه به جدول (۳)، تنها روش ARCO با دسته‌بند SVM و روش پیشنهادی در این مطالعه با دسته‌بند KNN دارای خطای صفر یا دقت ۱۰۰٪ هستند، با این تفاوت که با توجه به جدول (۴) تعداد ژن‌های انتخابی روش پیشنهادی با ۳۸ ژن کمتر از ۵۲ ژن روش ARCO است.

نکته مهم در مورد روش پیشنهادی نسبت به سایر روش‌ها، توانایی تفکیک ژن‌های موثر در هر کلاس است که می‌تواند کمک شایانی به فعالیت‌های آزمایشگاهی و داروسازی کند. جدول (۵) چگونگی توزیع ۳۸ ژن انتخابی در کلاس‌های مختلف سرطان پستان را نشان می‌دهد.

مناسب β ، به ماهیت داده‌ها و نوع مدل یادگیری استفاده شده برای دسته‌بندی ^{۲۰} بستگی دارد. برای پیاده‌سازی ایده‌های روش پیشنهادی شامل امتیازدهی و رتبه‌بندی نهایی ژن‌ها و استفاده از MATLAB دسته‌بندهای KNN و NB، از نرم‌افزار WEKA استفاده شده است. همچنین برای استفاده از دسته‌بند SVM از نرم‌افزار KNN استفاده شده است. زمان متوسط اجرای بخش اول شامل امتیازدهی به ژن‌ها در هر کلاس و رتبه‌دهی نهایی به آن‌ها روی سیستم با مشخصات سرعت پردازنده ۱/۱ گیگا هرتز و حافظه اصلی ۴ گیگا بايت با سیستم عامل ویندوز ۱۰، برابر ۱/۲۵ ثانیه است. همچنین زمان متوسط اجرای بخش دوم شامل استفاده از دسته‌بندهای NB و KNN و SVM به ترتیب برابر با ۳۶/۷۸ ثانیه، ۲/۲۷ ثانیه و ۲۳/۶۵ ثانیه است.

یافته‌ها

پس از محاسبه امتیاز هر ویژگی با استفاده از روابط (۱) و (۲) و تعیین ترتیب نهایی ویژگی‌ها در لیست خروجی، از مدل‌های یادگیری رایج و پرکاربرد SVM و KNN و NB، برای بررسی کارایی روش پیشنهادی بر اساس معیارهای دقت دسته‌بندی و تعداد ژن‌های انتخابی، استفاده شد و برای محدود کردن تعداد آزمایش‌ها برای انتخاب تعداد ویژگی‌های بهینه، حداقل تعداد ژن‌های انتخابی برابر ۱۰۰ در نظر گرفته شد. در آزمایش اول، مقادیر ژن ابتدایی لیست (ژن دارای رتبه اول) و در آزمایش دوم، مقادیر دو ژن ابتدایی لیست و در آزمایش kام $1 \leq k \leq 100$ ، مقادیر k ژن ابتدایی لیست از ۵۴ نمونه داده‌های آموزش برای آموزش هر یک از دسته‌بندها، استفاده می‌شود و هر بار برای بهدست آوردن دقت دسته‌بندی، مقادیر ژن‌های متناظر از ۳۰ نمونه‌ی مربوط به داده‌های آزمون به دسته‌بند داده می‌شود. و در پایان حداقل دقت بهدست آمده به همراه تعداد ژن‌های متناظر ثبت می‌شود. دقت یک روش دسته‌بندی، در صد نمونه‌های دسته‌بندی شده درست را در میان تمام نمونه‌ها نشان می‌دهد. نتایج بهدست آمده بر اساس مقادیر مختلف β در جدول‌های (۱) و (۲) نشان داده شده است.

^{۲۰} Classification

جدول ۳: حداقل میانگین خطای خطا در حداقل ۱۰۰ ژن انتخابی

FAST*	ARCO*	ReliefF*	Proposed method	دسته‌بند
۰/۱۴۴	۰/۰۱۴	۰/۱۶	۰/۶۶	NB
۰/۰۵۲	۰	۰/۱۴۳	۰/۰۴	SVM
۰/۰۸۱	۰/۰۱۹	۰/۰۷۴	۰	KNN

* نتایج گزارش شده در (۲۵)

جدول ۴: تعداد ژن‌های انتخابی هر روش با حداقل میانگین خطای خطا

FAST*	ARCO*	ReliefF*	Proposed method	دسته‌بند
۶۳	۶۸	۸۲	۱۷	NB
۷۹	۵۲	۵۷	۲۲	SVM
۶۲	۵۵	۹۷	۳۸	KNN

* نتایج گزارش شده در (۲۵)

جدول ۵: توزیع ۳۸ ژن انتخابی با استفاده از روش پیشنهادی در کلاس‌های مختلف سرطان پستان

کلاس سرطان								
شماره ژن								
G4450	G3150	G8506	G6797	G7961	G868	G7797	G9111	Luminal
	G2188	G1379	G6595	G1054	G7133	G886	G8921	ERBB2
G4805	G6243	G6410	G4211	G2452	G6623	G4380	G6802	Basal-like
G6554	G2951	G6606	G5342	G322	G1756	G74	G1280	Normal-like
	G7625	G2030	G8137	G5651	G6385	G8949	G4919	Cell-lines

بحث

در روش پیشنهادی در این مطالعه، تاثیر هر ژن در هریک از کلاس‌های سرطان پستان به طور جداگانه بررسی شد. امتیاز هر ژن در هر کلاس بر اساس چگونگی افزای اعضای آن کلاس توسط ژن مورد نظر، از طریق رابطه (۱) محاسبه شد. امتیاز به دست آمده برای در نظر گرفتن تاثیر ژن‌های مشترک، توسط رابطه (۲) اصلاح شد و در نهایت، ژن‌های موثر در هر کلاس به ترتیب امتیاز در لیست خروجی ژن‌های انتخابی قرار گرفتند. سپس از سه دسته‌بند SVM و KNN و NB، برای بررسی کیفیت ژن‌های انتخابی در تشخیص انواع سرطان پستان، استفاده شد. نتایج مدل‌سازی نشان داد که با استفاده از دسته‌بند KNN و فقط ۳۸ ژن از ۹۲۱۶ ژن اولیه، کلیه نمونه‌های سرطان پستان از مجموعه داده‌های آزمایش به درستی تشخیص داده شد. در این بخش کارایی روش پیشنهادی بر اساس معیارهای دقت دسته‌بندی و تعداد ژن‌های انتخابی، با تعدادی از روش‌های مرتبط که آن‌ها نیز روی مجموعه داده‌های مورد بحث این مقاله اعمال شده‌اند، مقایسه می‌شود.

همان‌گونه که در جدول (۵) مشخص است، روش پیشنهادی توانسته است ژن‌های موثر در هر کلاس را به صورت جداگانه شناسایی کند. علاوه براین، با بررسی بیشتر ژن‌های موثر در هر کلاس، ژن‌های با بیان بیشتر یا کمتر از حد نرمال در هر کلاس قابل تشخیص است. به عنوان مثال ژن G9111 در ۸۵٪ نمونه‌های آموزش مربوط به سرطان نوع Luminal یعنی ۱۷ نمونه از ۲۰ نمونه، بیان بیش از حد داشته است.

در مورد ژن G1054، در ۵ نمونه از ۶ نمونه مربوط به سرطان نوع ERBB2، بیان ژنی بیشتر نسبت به سایر نمونه‌ها مشاهده گردید. در سرطان نوع Basal-like، ۶ نمونه از ۷ نمونه داده‌های آموزش در ژن G2452، مقادیر بیان ژنی بالاتر نسبت به سایر نمونه‌های انواع دیگر سرطان داشتند. ژن‌های G1280 و G74 در کلیه نمونه‌های کلاس نوع Normal-like بیان بیش از حد نرمال داشته‌اند. در مورد کلاس Cell-lines، کلیه ژن‌های مشخص شده، بیان ژن کمتر از حد نرمال در ۱۲ نمونه مربوطه داشتند.

ژو و همکاران در (۳۰)، با بهره‌گیری از مفهوم آنتروپی و مبانی تئوری اطلاعات، روشی به نام IK-TSP را که توسعه یافته روش K-TSP^{۲۶} برای استفاده در مسائل دسته‌بندی چند کلاسه است، ارایه کردند که با ۵۶ زن انتخابی، انواع سرطان سینه را با دقت ۸۳/۳۳٪ پیش‌بینی کردند.

چن^{۲۹} و همکاران در (۳۱)، روشی با نام^{۳۰} RS-based DC ارایه دادند که در آن ابتدا مسئله دسته‌بندی چند کلاسه را بر اساس روش OVR^{۳۱}، به تعدادی مسئله بایتری تبدیل کردند. سپس معیارهایی بر اساس مبانی تئوری اطلاعات و آنتروپی برای ارزیابی افقی و عمودی مقادیر یک زن نسبت به سایر زن‌ها ارایه دادند که بر اساس آن‌ها در کنار معیاری برای ارزیابی نسبی مقادیر درون زنی، امتیاز زن‌ها را تعیین کردند. چن و همکاران با این روش با تعداد ۱۵ زن انتخابی به دقت ۹۳/۳۳٪ در تشخیص درست سرطان پستان دست یافتند که از نظر تعداد زن انتخابی نسبت به روش پیشنهادی عملکرد بهتری دارد ولی در مقایسه با دقت ۱۰۰ درصدی روش پیشنهادی، کارایی پایین‌تری دارد.

ونگ^{۳۲} و وی^{۳۳} (۳۲) در سال ۲۰۱۷ دو روش انتخاب ویژگی DSCFS و CAM را بر اساس اندازه‌گیری توانایی یک ویژگی در دسته‌بندی زیر مسایل دو کلاسه و مفهوم مکمل بودن ویژگی‌ها ارایه کردند. آن‌ها با ۶ زن انتخاب شده توسط روش DSCFS و با استفاده از سه دسته‌بند NB و SVM و KNN به ترتیب به دقت‌های ۹۰ و ۹۶/۶۷ و ۹۶/۶۷ درصد دست یافتند. آن‌ها همچنین با ۱۵۲ زن انتخاب شده توسط روش CAM و با استفاده از سه دسته‌بند NB و SVM و KNN به ترتیب به دقت‌های ۸۰ و ۹۶/۶۷ و ۹۰ درصد را گزارش کردند. با وجود برتری روش DSCFS در تعداد زن انتخابی، روش پیشنهادی با توجه به معیار دقت، عملکرد بهتری دارد.

سان^{۳۴} و همکاران (۳۳) در سال ۲۰۱۹ روشی به نام ECOC-MDC برای تبدیل مسئله دسته‌بندی چند کلاسه به دسته‌بندی‌های دوکلاسه بر اساس مفاهیم

در سال ۲۰۰۵، Tan^{۲۱} و همکاران در (۲۶) روش‌های HC-K-TSP و HC-TSP مقادیر جفت زن‌ها ارایه کردند که توسعه یافته روش TSP^{۲۲} (۲۷) برای استفاده در دسته‌بندی‌های چند کلاسه است. آن‌ها با استفاده از روش‌های فوق به ترتیب با انتخاب ۲۴ زن، به دقت ۶۶/۶۷٪ دست یافتند که هرچند از نظر تعداد زن عملکرد بهتری را نسبت به ۳۸ زن انتخابی روش پیشنهادی نشان می‌دهد ولی از نظر دقت، فاصله زیادی تا دقت ۱۰۰٪ دارد. Tan و همکاران همچنین با نرم‌افزار PAM^{۲۳} (۲۸) که یک ابزار تحلیل آماری رایج برای کار روی داده‌های ریزآرایه است، با ۴۸۲۲ زن انتخابی، انواع سرطان سینه را با دقت ۹۳/۳۳٪ درصد پیش‌بینی کردند که با وجود بهبود دقت دسته‌بندی، با تعداد ۴۸۲۲ زن انتخابی عملکرد مناسبی ندارد.

ونگ^{۲۴} و همکاران در (۲۹)، الگوریتم^{۲۵} TSG را ارایه دادند که بر اساس انتخاب مجموعه‌هایی از k زوج ویژگی با امتیاز بالاتر براساس معیار chi-square و ارزیابی هر یک از مجموعه ویژگی‌ها بر پایه‌ی LOOCV^{۲۶}، به دقت ۸۶/۶۷٪ با تعداد ۶۳ زن دست یافتند که در هر دو پارامتر دقت دسته‌بندی و تعداد زن انتخابی، عملکرد ضعیف‌تری نسبت به روش پیشنهادی دارد.

لی سان^{۲۷} و همکاران در (۲۵) در سال ۲۰۱۷ یک روش ROC^{۲۸} انتخاب ویژگی به نام AVC را بر اساس منحنی ROC و خاصیت مکمل‌پذیری ارایه کردند. آن‌ها از چهار دسته‌بند SVM و C4.5 و KNN و NB برای ارزیابی ۷۲ زن انتخاب شده توسط AVC استفاده کردند. نتایج گزارش شده توسط آن‌ها با استفاده از چهار دسته‌بند ذکر شده، به ترتیب خطای میانگین ۰/۰۰۶ و ۰/۱۳۶ و ۰/۰۱۹ و ۰/۰۷۳ درصد را در تشخیص انواع سرطان پستان نشان داد که با وجود دستیابی به دقت نزدیک ۱۰۰ درصد، تعداد ۷۲ زن انتخابی بیشتر از ۳۸ زن روش پیشنهادی است.

²¹ Tan

²² Top scoring pair

²³ Prediction analysis of microarrays

²⁴ Wang

²⁵ Top scoring genes

²⁶ Leave-one-out cross validation

²⁷ Wang

²⁸ Receiver Operator Characteristic

²⁹ Chen

³⁰ Relative simplicity based direct classification

³¹ One versus rest

³² Wang

³³ Wei

³⁴ sun

در این مطالعه، از ارتباط بین بیان غیرنرم‌مال ژن و سرطان به عنوان دانش در مسئله انتخاب ژن در داده‌های ریزآرایه استفاده شد که باعث گردید روش پیشنهادی علاوه بر تشخیص ۱۰۰ درصدی انواع سرطان پستان، قادر به تشخیص ژن‌های انتخابی باشد. توانایی تفکیک ژن‌ها در هر کلاس سرطان به علاوه تشخیص بیان بیشتر از حد یا کمتر از حد ژن‌های انتخابی، از خصوصیات ویژه روش پیشنهادی است که می‌تواند مورد استفاده متخصصین و پژوهشگران حوزه تشخیص و درمان بیماری سرطان پستان قرار گیرد.

تقدیر و تشکر

این پژوهش نتیجه بخشی از فعالیت‌های صورت گرفته در قالب رساله دکتری با عنوان "انتخاب ویژگی مبتنی بر تئوری اطلاعات برای انتخاب ژن‌های موثر در داده‌های ریزآرایه" است که در تاریخ ۱۳۹۶/۱۱/۱۱ در شورای تحصیلات تکمیلی دانشگاه یزد با شماره ۳۹۶۹۰ به تصویب رسیده است. از کلیه همکاران گرامی آزمایشگاه هوش محاسباتی دانشگاه یزد که در انجام این پژوهش همراهی نمودند، تشکر و قدردانی می‌شود.

تعارض منافع

نویسنده‌گان اعلام می‌دارند که هیچ تعارض منافعی در پژوهش حاضر وجود ندارد.

پیچیدگی داده‌ها ارایه دادند. آن‌ها از سه روش انتخاب Wilcoxon و ROC و T-test و ویژگی ROC با ارزش‌تر استفاده کردند. آن‌ها با استفاده از روش ECOC-MDC و دسته‌بند NB و با ۱۰۰ ژن انتخابی الگوریتم ROC، با دقت ۹۶٪ انواع سرطان پستان را تشخیص دادند. همچنین با دسته‌بند SVM و با ۱۰۰ ژن انتخابی الگوریتم T-test، به دقت ۹۷/۳٪ و با دسته‌بند SVM و با ۱۰۰ ژن انتخابی الگوریتم Wilcoxon، به دقت ۹۸/۷٪ دست یافتند.

در مجموع دست یافتن به دقت ۱۰۰ درصدی، نقطه قوت روش پیشنهادی محسوب می‌شود. هر چند از نظر تعداد ژن‌های انتخابی، برخی روش‌های ذکر شده عملکرد بهتری داشتند ولی با در نظر گرفتن هم‌زمان دو پارامتر دقت دسته‌بندی و تعداد ژن انتخابی، روش پیشنهادی عملکرد قابل قبولی نسبت به سایر روش‌ها دارد.

علاوه بر این، با توجه به این‌که امتیاز هر ژن در هر کلاس به صورت مجزا تعیین می‌شود، می‌توان در صورت نیاز از این روش برای یافتن ژن‌های مشترک بین کلاس‌های بیماری استفاده کرد. به این ترتیب که ژن‌هایی که دارای امتیاز بالا در دو یا چند کلاس باشند، می‌توانند به عنوان ژن‌های نامزد برای مشخص شدن ژن‌های موثر مشترک، مورد مطالعه قرار گیرند.

نتیجه‌گیری

References

- Smith RA, Andrews KS, Brooks D, Fedewa SA, Manassaram-baptiste D, Saslow D, et al. Cancer Screening in the United States, 2018: A Review of Current American Cancer Society Guidelines and Current Issues in Cancer Screening. CA Cancer J Clin. 2018;68(4):297-316.
- Enayat R, Salehiniya H. An investigation of changing patterns in breast cancer incidence trends among Iranian women. J Sabzevar Univ Med Sci. 2015;22(1):27-35.
- Wang YA, Johnson SK, Brown BL, McCarragher LM, Al-Sakkaf K, Royds JA, et al. Enhanced anti-cancer effect of a phosphatidylinositol-3 kinase inhibitor and doxorubicin on human breast epithelial cell lines with different p53 and oestrogen receptor status. Int J Cancer. 2008;123(7):1536-44.
- Sheikhpour R, Ghassemi N, Yaghmaei P, Ardekani JM, Shiryazd M. Immunohistochemical assessment of P53 protein and its correlation with clinicopathological characteristics in breast cancer patients. Indian J Sci Technol. 2014;7(4):472-9.
- Vallejos CS, Gómez HL, Cruz WR, Pinto JA, Dyer RR, Velarde R, et al. Breast Cancer Classification According to Immunohistochemistry Markers: Subtypes and Association With Clinicopathologic Variables in a Peruvian Hospital Database. Clin Breast Cancer. 2010;10(4):294-300.
- Cheang MCU, Chia SK, Voduc D, Gao D, Leung S, Snider J, et al. Ki67 index, HER2 status, and

- prognosis of patients with luminal B breast cancer. *JNCI J Natl Cancer Inst.* 2009; 101(10): 736-50.
7. Perou CM, Sørlie T, Eisen MB, Van De Rijn M, Jeffrey SS, Rees CA, et al. Molecular portraits of human breast tumours. *Nature.* 2000; 406(6797): 747-52.
 8. Sørlie T, Perou CM, Tibshirani R, Aas T, Geisler S, Johnsen H, et al. Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications. *Proc Natl Acad Sci.* 2001; 98(19):10869-74.
 9. Dai X, Li T, Bai Z, Yang Y, Liu X, Zhan J, et al. Breast cancer intrinsic subtype classification, clinical use and future trends. *Am J Cancer Res.* 2015; 5(10):2929-43.
 10. Piatetsky-Shapiro G, Tamayo P. Microarray data mining: facing the challenges. *ACM SIGKDD Explor Newslet.* 2003; 5(2):1-5.
 11. Zhang L, Lin X. Some considerations of classification for high dimension low-sample size data. *Stat Methods Med Res.* 2013; 22(5):537-50.
 12. Guyon I, Elisseeff A. An Introduction to Variable and Feature Selection. *J Mach Learn Res.* 2003; 3(3):1157-82.
 13. Glazier AM. Finding Genes That Underlie Complex Traits. *Science.* 2002; 298(5602): 2345-9.
 14. Saeys Y, Inza I, Larrañaga P. A review of feature selection techniques in bioinformatics. *Bioinformatics.* 2007; 23(19):2507-17.
 15. Golub TR, Slonim DK, Tamayo P, Huard C, Gaasenbeek M, Mesirov JP, et al. Molecular classification of cancer: Class discovery and class prediction by gene expression monitoring. *Science.* 1999; 286(5439):531-7.
 16. Tsafrir D, Bacolod M, Selvanayagam Z, Tsafrir I, Shia J, Zeng Z, et al. Relationship of gene expression and chromosomal abnormalities in colorectal cancer. *Cancer Res.* 2006; 66(4): 2129-37.
 17. Prelich G. Gene overexpression: uses, mechanisms, and interpretation. *Genetics.* 2012; 190(3): 841-54.
 18. Gray JW, Collins C. Genome changes and gene expression in human solid tumors. *Carcinogenesis.* 2000; 21(3):443-52.
 19. Martincorena I, Campbell PJ. Somatic mutation in cancer and normal cells. *Science (80).* 2015; 349(6255):1483-9.
 20. Zainulabdeen A, Yao P, Zare H. Underexpression of specific interferon genes is associated with poor prognosis of melanoma. *PLoS One.* 2017; 12(1):e0170025.
 21. Pinter CC. *A Book of SET THEORY.* Dover Publications; 2014.
 22. Zhao Z, Wang L, Liu H, Ye J. On Similarity Preserving Feature Selection. *IEEE Trans Knowl Data Eng.* 2013; 25(3):619-32.
 23. Wang R, Tang K. Feature Selection for Maximizing the Area Under the ROC Curve. In: *2009 IEEE International Conference on Data Mining Workshops.* 2009; 400-5.
 24. Chen X, Wasikowski M. FAST: A Roc-based Feature Selection Metric for Small Samples and Imbalanced Data Classification Problems. In: *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.* 2008; 124-32.
 25. Sun L, Wang J, Wei J. AVC: Selecting discriminative features on basis of AUC by maximizing variable complementarity. *BMC bioinformatics.* 2017; 18(3):50.
 26. Tan AC, Naiman DQ, Xu L, Winslow RL, Geman D. Simple decision rules for classifying human cancers from gene expression profiles. *Bioinformatics.* 2005; 21(20):3896-904.
 27. Geman D, d'Avignon C, Naiman DQ, Winslow RL. Classifying gene expression profiles from pairwise mRNA comparisons. *Stat Appl Genet Mol Biol.* 2004; 3(1):1-19.
 28. Tibshirani R, Hastie T, Narasimhan B, Chu G. Diagnosis of multiple cancer types by shrunken centroids of gene expression. *Proc Natl Acad Sci.* 2002; 99(10):6567-72.
 29. Wang H, Zhang H, Dai Z, Chen M, Yuan Z. TSG: a new algorithm for binary and multi-class cancer classification and informative genes selection. *BMC Med Genomics.* 2013; 6(1):S3.
 30. Zhou C, Wang S, Blanzieri E, Liang Y. An entropy-based improved k-top scoring pairs (TSP) method for classifying human cancers. *African J Biotechnol.* 2012;11(45):10438-45.
 31. Chen Y, Wang L, Li L, Zhang H, Yuan Z. Informative gene selection and the direct classification of tumors based on relative simplicity. *BMC Bioinformatics.* 2016;17(1):44.
 32. Wang S, Wei J. Feature selection based on measurement of ability to classify subproblems. *Neurocomputing.* 2017; 224:155-65.
 33. Sun M, Liu K, Wu Q, Hong Q, Wang B, Zhang H. A novel ECOC algorithm for multiclass microarray data classification based on data complexity analysis. *Pattern Recognit.* 2019; 90: 346-62.