

Two Stages Sparse Representation-based Classifier and its Application for Cancer Classification

M. Miri¹, M. T. Sadeghi^{2*}, V. Abootalebi³

¹M.Sc., Signal Processing Research Lab, Electrical and Computer Engineering Department, Yazd University, Yazd, Iran, m.miri@stu.yazd.ac.ir

^{2*} Assistant Professor, Signal processing Research Lab, Electrical and Computer Engineering Department, Yazd University, Yazd, Iran.

³ Assistant Professor, Signal processing Research Lab, Electrical and Computer Engineering Department, Yazd University, Yazd, Iran, abootalebi@yazd.ac.ir

Abstract

Successful outcomes of Sparse Representation-based Classifier (SRC) and Sparse Subspace Clustering (SSC) in many applications motivated us to combine these methods and propose a hierarchical classifier. The main idea behind the SRC and SSC algorithms is to represent a data using a sparse linear combination of elementary signals so that those elementary signals which are similar to the data contribute mainly in the representation. In this paper, the performance of a sparse representation based classifier is improved by pre-clustering of training samples using the SSC algorithm. A two-stage SRC is then designed using the resulting clusters. A test data is classified by first determining the most similar cluster. The data label is subsequently found using the second stage classifier. The performance of the proposed method is evaluated considering cancer classification problem using the 14-Tumors microarray dataset. Due to low number of data samples per each class and high dimensionality of the data, this is a challenging problem. Curse of dimensionality, overfitting of the classifier to the training data and computational complexity are the possible related problems. Our experimental results show that the proposed method outperforms some other state of the art classifiers.

Key words: Sparse Subspace Clustering, Microarray data, Cancer classification, Hierarchical classifier, Sparse Representation-based Classification, Sparse representation.

* Corresponding author

Address: Electrical and Computer Engineering Department, Yazd University, P.O.Box: 89195-741, Postal Code: 77141, Yazd, I.R. Iran

Tel: +98 351 8122389

Fax: +98 351 8200144

E-mail: m.sadeghi@yazd.ac.ir

طبقه‌بندی‌کننده دو مرحله‌ای مبتنی بر نمایش تنک و کاربرد آن در تشخیص سرطان

ملیحه میری^۱، محمدتقی صادقی^{۲*}، وحید ابوطالبی^۳

^۱ دانشجوی کارشناسی ارشد، گروه مخابرات، دانشکده مهندسی برق و کامپیوتر، دانشگاه یزد، یزد، ایران m.miri@stu.yazd.ac.ir

^۲ استادیار، گروه مخابرات، دانشکده مهندسی برق و کامپیوتر، دانشگاه یزد، یزد، ایران.

^۳ استادیار، گروه مخابرات، دانشکده مهندسی برق و کامپیوتر، دانشگاه یزد، یزد، ایران abootalebi@yazd.ac.ir

چکیده

با توجه به نتایج موفقیت‌آمیز طبقه‌بندی‌کننده مبتنی بر نمایش تنک (SRC) و خوشه‌بندی زیرفضای تنک (SSC) در کاربردهای مختلف، در این مقاله با ترکیب این دو روش، یک روش طبقه‌بندی سلسله مراتبی ارائه می‌شود. ایده اصلی در روش‌های طبقه‌بندی و خوشه‌بندی مبتنی بر نمایش تنک، نمایش هر داده به صورت ترکیب خطی تنک از سایر داده‌ها است به گونه‌ای که داده‌های مشابه با داده مورد نظر در این ترکیب خطی بیشترین وزن را به خود اختصاص دهند. در روش پیشنهادی، به منظور دستیابی به صحت طبقه‌بندی بیشتر، ابتدا داده‌های آموزشی با استفاده از روش خوشه‌بندی زیرفضای تنک بخش‌بندی می‌شوند. سپس با استفاده از شیوه بکار گرفته شده در طبقه‌بندی‌کننده مبتنی بر نمایش تنک، طبقه‌بندی‌کننده‌ای دو مرحله‌ای طراحی می‌شود. در مرحله اول، خوشه‌ای که داده ورودی بیشترین شباهت را با آن دارد تعیین شده و در مرحله بعد طبقه مربوطه (برچسب داده) تعیین می‌شود. برای ارزیابی روش پیشنهادی از دادگان ریزآرایه 14-Tumors - که حاوی اطلاعات مربوط به ۱۴ نوع سرطان مختلف است - استفاده شده است. از جمله ویژگی‌های این دادگان تعداد زیاد بعد در مقابل تعداد کم نمونه در هر دسته است که عمل طبقه‌بندی آن‌ها را به مسأله‌ای چالش‌برانگیز تبدیل می‌کند. ابعاد زیاد داده‌ها نه تنها مشکلاتی از جمله نفرین ابعاد و بیش انطباق طبقه‌بندی‌کننده به داده‌های آموزشی را به دنبال دارد، بلکه باعث افزایش پیچیدگی محاسباتی شده؛ زمان لازم را برای اجرای الگوریتم‌ها افزایش می‌دهد. آزمایش‌های انجام شده بر این دادگان با استفاده از روش پیشنهادی نشان می‌دهد که در مقایسه با سایر روش‌های طبقه‌بندی، این روش به نتایج بهتری منجر می‌شود.

کلیدواژه‌ها: خوشه‌بندی زیرفضای تنک، دادگان ریزآرایه، طبقه‌بندی سرطان، طبقه‌بندی‌کننده سلسله مراتبی، طبقه‌بندی‌کننده مبتنی بر نمایش تنک، نمایش تنک.

*عهده‌دار مکاتبات

نشانی: یزد، صفائیه، چهارراه پژوهش، دانشگاه یزد، دانشکده مهندسی برق و کامپیوتر، کد پستی: ۷۷۱۴۱

تلفن: ۰۳۵۱-۸۱۲۲۳۸۹، دورنگار: ۰۳۵۱-۸۲۰۰۱۴۴، پیام نگار: m.sadeghi@yazd.ac.ir

۱- مقدمه

در سال‌های اخیر شاهد توسعه روزافزون استفاده از مفاهیم نمونه‌برداری فشرده (CS)^۱ و نمایش تنک (SR)^۲ در کاربردهای مختلف پردازش سیگنال هستیم [۲،۱] و تحقیقات بسیاری در این زمینه شکل گرفته است. در بسیاری از کاربردهای شناسایی الگو و بینایی ماشین با داده‌هایی با ابعاد زیاد از قبیل تصاویر چهره و یا ریزآرایه‌ها سروکار داریم.

بعد زیاد داده‌ها نه تنها موجب افزایش حافظه مورد نیاز و پیچیدگی محاسباتی الگوریتم‌ها می‌شود، به دلیل اثر نویز و تعداد کم داده‌ها در مقایسه با تعداد ابعاد، موجب کاهش کارایی الگوریتم‌ها نیز می‌شود که از این موضوع به عنوان نفرین ابعاد^۳ یاد می‌شود [۳]. داده‌های با ابعاد زیاد اغلب در زیرفضاهایی با بعد کمتر گسترده شده‌اند. روش‌های مختلف خوشه‌بندی^۴ به عنوان ابزاری برای جداسازی داده‌های متناظر با هر یک از این زیرفضاها به کار می‌روند.

اخیراً به استفاده از روش‌های مبتنی بر نمایش تنک در خوشه‌بندی داده‌ها نیز توجه بسیار شده است [۵،۴]. از بین تمام نمایش‌های ممکن برای یک داده بر حسب سایر داده‌ها، یک نمایش تنک متناظر با انتخاب تعداد کمی از داده‌ها از همان زیرفضای داده مربوطه است. چنین نمایشی با حل مسأله بهینه‌سازی تنک بدست می‌آید. روش‌های متداول خوشه‌بندی از جمله روش‌های تکراری [۶]، روش‌های جبری [۷]، MPPCA^۵ [۸]، RANSAC^۶ [۹] و روش‌های مبتنی بر خوشه‌بندی طیفی^۷ [۱۰] دارای محدودیت‌هایی هستند که استفاده از آنها را در عمل با مشکل مواجه می‌کند. برخی از این روش‌ها نیاز به داشتن اطلاعات اولیه‌ای از جمله تعداد و بعد خوشه‌ها دارند. همچنین بار محاسباتی زیاد - که بصورت نمایی با تعداد و بعد خوشه‌ها افزایش می‌یابد- و عدم مقاومت نسبت به نویز و خطای مدل‌سازی از دیگر معایب این روش‌ها است [۱۱]. در مقایسه با این روش‌ها، روش خوشه‌بندی زیرفضای تنک (SSC)^۸ که از حل معادلات بهینه‌سازی محدب^۹ برای یافتن نمایش‌های تنک استفاده

می‌کند، از پیچیدگی محاسباتی کمتری برخوردار است. همچنین این روش به مقداردهی اولیه نیاز ندارد و نسبت به نویز نیز مقاوم است [۴].

در سال ۲۰۰۹، رایت^{۱۱} استفاده از نمایش تنک سیگنال‌ها در طبقه‌بندی را پیشنهاد کرد، این روش - که طبقه‌بندی مبتنی بر نمایش تنک (SRC)^{۱۱} نامیده می‌شود- به نتایج رضایت‌بخشی بویژه در سیستم‌های تشخیص چهره منجر شده است [۱۲]. ایده اصلی این روش آن است که هر تصویر چهره به صورت ترکیب خطی تنک از سایر تصاویر چهره قابل بیان است و بیشترین ضرایب این ترکیب خطی متعلق به تصاویری است که در دسته‌ای یکسان با تصویر ورودی قرار دارند؛ لذا با صفر کردن سایر ضرایب و تنها استفاده از این ضرایب بزرگتر می‌توان به بازسازی قابل قبولی از تصویر ورودی دست یافت. با استفاده از معیار فاصله اقلیدسی میزان شباهت تصویر بازسازی شده با کمک تصاویر دسته‌های مختلف، با تصویر ورودی تعیین می‌شود و برچسب دسته‌ای که کمترین خطای بازسازی را تولید کند، به تصویر ورودی تعلق می‌گیرد. پیشرفت‌های اخیر در زمینه زیست‌شناسی مولکولی و فناوری ریزآرایه^{۱۲}، امکان پایش سطح بیان^{۱۳} هزاران ژن (میزان فعال بودن ژن‌ها) را به صورت همزمان فراهم کرده است و مقادیر هنگفتی را از داده‌های ریزآرایه تولید کرده است. این داده‌ها در تشخیص و طبقه‌بندی انواع بافت‌های سرطانی نقش به‌سزایی دارند؛ اما مهمترین چالش در مورد این دادگان را می‌توان ابعاد بسیار بالای آنها در مقایسه با تعداد کم نمونه‌ها دانست که طراحی طبقه‌بندی‌کننده‌های مناسب را دچار مشکل می‌کند. اگرچه تاکنون تلاش‌های بسیاری با هدف یافتن طبقه‌بندی‌کننده‌های با صحت زیاد و پیچیدگی محاسباتی کم انجام شده است، بیشتر مطالعاتی که در این زمینه انجام شده توجه خود را به دادگانی با تعداد دسته‌های کم (دو یا سه نوع سرطان مختلف) معطوف کرده‌اند و از روش‌هایی مثل تحلیل جداکننده خطی (LDA)^{۱۴}، شبکه‌های عصبی، خوشه‌بندی، نزدیک‌ترین همسایگی (NN)^{۱۵} و ماشین بردار پشتیبان (SVM)^{۱۶} استفاده کرده‌اند [۱۵،۱۴،۱۳]. دادگان 14-Tumors -

^۱Compressive Sensing^۲Sparse Representation^۳Curse of Dimensionality^۴Clustering^۵Mixture of Probabilistic Principal Component Analysis^۶Random Sample Consensus^۷Spectral Clustering^۸Sparse Subspace Clustering^۹Convex Optimization^{۱۰}Wright^{۱۱}Sparse Representation-based Classification^{۱۲}Microarray^{۱۳}Expression^{۱۴}Linear Discriminant Analysis^{۱۵}Nearest Neighbor^{۱۶}Support Vector Machine

می‌شود و در مرحله بعد با توجه به خوشه تعیین شده، عمل طبقه‌بندی انجام شده؛ برچسب داده آزمایشی تعیین می‌شود. با توجه به این که یکی از روش‌های متداول در دسته‌بندی انواع سرطان‌ها خوشه‌بندی داده‌های مربوط به آنها است، انتظار می‌رود روش پیشنهادی - که در مرحله اول داده‌های آموزش را خوشه‌بندی و سپس آن را طبقه‌بندی می‌کند - به نتایج مناسبی در این حوزه دست یابد. مقایسه نتایج به دست آمده از روش معرفی شده با نتایج حاصل از تحقیقات قبلی نشان دهنده عملکرد قابل قبول سیستم پیشنهادی است.

ساختار این مقاله به قرار زیر است: در بخش ۲ مقدمه‌ای بر نمایش تنک بیان شده است، سپس به توضیح اجمالی طبقه‌بندی‌کننده مبتنی بر نمایش تنک و خوشه‌بندی زیرفضای تنک پرداخته شده است. در بخش ۳ روش پیشنهادی و در بخش ۴ توضیح مختصری در مورد پایگاه داده استفاده شده و همچنین نتایج آزمایش‌های انجام شده آورده شده است. در نهایت در بخش ۵ جمع‌بندی و نتیجه‌گیری ارائه شده است.

۲- نمایش تنک (SR)

امروزه نمایش تنک به ابزاری قوی برای نمایش و فشرده‌سازی سیگنال‌ها تبدیل شده است. بر اساس این نمایش می‌توان سیگنالی را بصورت ترکیب خطی تنک از داده‌های آموزشی تقریب زد [۲]. در حالت کلی ترکیبی خطی از داده‌های آموزشی را می‌توان به صورت یک دستگاه معادلات خطی به شکل $\mathbf{y} = \mathbf{Ys}$ نشان داد. که در آن ماتریس $\mathbf{Y} = [\mathbf{Y}_1, \dots, \mathbf{Y}_c]$ - که از کنار هم قرار گرفتن تمام داده‌های آموزشی مربوط به \mathbf{c} دسته تشکیل شده است - ماتریس دیکشنری نامیده شده و ستون‌های آن اتم نامیده می‌شوند. اگر در این دستگاه تعداد معادلات از تعداد مجهولات کمتر باشد (تعداد ستون‌های \mathbf{Y} بیشتر از تعداد سطرهای آن باشد)، دستگاه از نوع فرومعین^{۲۲} بوده و دارای بی‌شمار جواب است. با اعمال شرط تنکی بر جواب حاصل از دستگاه فوق می‌توان به پاسخی یکتا دست یافت. از آنجا که نرم صفر یک بردار معرف تعداد مؤلفه‌های غیر صفر آن است، کمینه‌سازی نرم صفر منجر به تحمیل شرط تنکی بر بردار خواهد شد؛ بنابراین

که در این تحقیق بررسی شده‌اند - شامل اطلاعات مربوط به چهارده نوع سرطان مختلف است که تا کنون اغلب از دید انتخاب ویژگی^{۱۷} بررسی شده است [۱۶]. از جمله پژوهش‌های انجام شده در طبقه‌بندی این پایگاه داده می‌توان به روش ارائه شده در مرجع [۱۶] اشاره کرد. در این مقاله مسأله طبقه‌بندی چهارده دسته سرطان مختلف به چند مسأله طبقه‌بندی باینری تبدیل، و از طبقه‌بندی‌کننده‌های مختلفی نیز استفاده شده است که بهترین عملکرد مربوط به طبقه‌بندی‌کننده SVM بوده است. در طبقه‌بندی باینری با روش‌های ^{18}OVO و ^{19}OVA عمل طبقه‌بندی انجام شده است که به ترتیب از $\frac{c(c-1)}{2}$ و c ابرصفحه برای جداسازی دسته‌ها استفاده می‌کنند که c تعداد دسته‌ها است. در منابع [۱۷، ۱۸] با استفاده از روش‌های مبتنی بر شبکه عصبی چهارده نوع سرطان طبقه‌بندی شده‌اند؛ به عنوان مثال در مرجع [۱۷] از دو طبقه‌بندی‌کننده ^{20}ANN متوالی برای طبقه‌بندی استفاده شده است. در این روش ابتدا با استفاده از طبقه‌بندی‌کننده اول، دو دسته‌ای که بیشترین احتمال تعلق داده آموزش به آن‌ها وجود دارد انتخاب شده؛ سپس با اعمال طبقه‌بندی‌کننده دوم و انجام طبقه‌بندی باینری بین دو دسته انتخابی، برچسب داده آموزش تعیین می‌شود.

طبقه‌بندی‌کننده‌ای فازی نیز در مطالعه [۱۹] ارائه شده است (CD-MFS)^{۲۱}، که با استفاده از قوانین فازی، طبقه‌بندی را انجام می‌دهد. مرجع [۲۰] نیز از روشی مبتنی بر نمایش تنک با عنوان SR برای جداسازی دسته‌های مختلف استفاده کرده است و به نتایج بهتری در مقایسه با تعمیم‌های مختلف روش SVM دست یافته است.

با وجود تحقیقاتی که در این زمینه شکل گرفته‌اند، همچنان دسته‌بندی این دادگان با مقادیر خطای قابل توجهی همراه است و به صورت بالینی قابل استفاده نیست.

در این مقاله با ترکیب دو روش خوشه‌بندی و طبقه‌بندی مبتنی بر نمایش تنک، روش طبقه‌بندی سلسله مراتبی ارائه شده است. بدین منظور در مرحله آموزش، داده‌های آموزشی خوشه‌بندی شده؛ سپس آزمایش سیستم در دو مرحله انجام می‌شود. در مرحله اول خوشه مربوط به داده ورودی مشخص

¹⁷Compressive Sensing¹⁸One-Versus-One¹⁹One-Versus-All²⁰Artificial Neural Networks²¹Cancer Diagnosis with Memetic Fuzzy System²²Under-Determined

می‌شوند. در این روش در هر مرحله یک اتم از دیکشنری که بیشترین شباهت را به داده آزمون دارد، به عنوان عضو فعال در ترکیب خطی در نظر گرفته شده؛ ضریب مربوط به آن محاسبه می‌شود. تفاضل حاصل ضرب این تقریب ۱- تنک در دیکشنری از داده آزمون را به عنوان باقی‌مانده در نظر گرفته؛ مراحل فوق تکرار می‌شوند. در هر مرحله جمع تقریب‌های ۱- تنک به دست آمده با تقریب‌های قبلی به عنوان تقریب جدید در نظر گرفته می‌شوند و این روند تا جایی ادامه می‌یابد که یا تعداد مراحل مشخصی طی شود و یا خطا از مقدار معینی کمتر شود. در روش OMP^{۲۷} - که تعمیم‌یافته روش MP است- در هر مرحله ضرایب ستون‌های فعال از ماتریس دیکشنری به صورت مستقل از نتایج مراحل قبل انتخاب می‌شوند و از نتایج قبلی تنها در یافتن مکان مؤلفه‌های غیرصفر استفاده می‌شود. از جمله مهم‌ترین مزایای این روش‌ها می‌توان به سرعت زیاد آن‌ها اشاره کرد [۲۳]. از جمله دیگر الگوریتم‌هایی که در محاسبه نمایش تنک استفاده می‌شوند، می‌توان به Homotopy، PALM^{۲۸} و DALM^{۲۹} اشاره کرد [۲۴].

اگر N_i داده آموزشی متعلق به دسته i ام در ستون‌های ماتریس $\mathbf{Y}_i \in \mathbf{R}^{D \times N_i}$ به صورت $\mathbf{Y}_i = [\mathbf{y}_{i,1}, \mathbf{y}_{i,2}, \dots, \mathbf{y}_{i,N_i}]$ قرار گرفته باشند؛ بر اساس نمایش تنک چنانچه $\mathbf{y} \in \mathbf{R}^D$ داده ورودی متعلق به دسته i ام باشد، می‌توان آن را برحسب ستون‌های ماتریس \mathbf{Y}_i نمایش داد:

$$\begin{aligned} \mathbf{y} &= s_{i,1}\mathbf{y}_{i,1} + s_{i,2}\mathbf{y}_{i,2} + \dots + s_{i,N_i}\mathbf{y}_{i,N_i} \\ &= \sum_{j=1}^{N_i} s_{i,j}\mathbf{y}_{i,j} \end{aligned} \quad (3)$$

که در این رابطه $s_{i,j}$ ها ضرایب اسکالر هستند. از این ایده در طبقه‌بندی‌کننده مبتنی بر نمایش تنک استفاده می‌شود.

۳-۲-۱- طبقه‌بندی مبتنی بر نمایش تنک (SRC)
طبقه‌بندی با استفاده از نمایش تنک سیگنال شامل دو مرحله است: ابتدا داده ورودی با استفاده از ماتریس دیکشنری بصورت تنک کد می‌شود؛ سپس با استفاده از این ضرایب تنک و با بازسازی داده آزمایشی، طبقه‌بندی انجام

از رابطه زیر برای به دست آوردن نمایش تنک یک بردار بر نمونه‌های آموزشی استفاده می‌شود:

$$\min \|\mathbf{s}\|_0 \quad s.t. \quad \mathbf{y} = \mathbf{Y}\mathbf{s} \quad (1)$$

\mathbf{y} سیگنال ورودی و \mathbf{s} بردار ضرایب ترکیب خطی است و $\|\mathbf{s}\|_0$ نیز تعداد مؤلفه‌های غیر صفر آن را نشان می‌دهد. \mathbf{s} را \mathbf{k} -تنک گویند اگر دارای حداکثر \mathbf{k} مؤلفه غیرصفر باشد. همان‌گونه که قبلاً ذکر شد اگر ماتریس دیکشنری فراکامل^{۳۳} باشد، دستگاه معادلات فوق پاسخی یکتا و تنک خواهد داشت. چون نرم صفر محدب و مشتق‌پذیر نیست، مسأله فوق NP-Complete^{۳۴} بوده؛ حل آن منوط به جستجوی کامل است که در ابعاد زیاد غیر ممکن است. روشی رایج در حل سیستم‌های خطی با فرض تنک بودن جواب، تقریب نرم صفر با نرمی از مرتبه بالاتر است که کمینه کردن آن ساده‌تر باشد. بدین منظور برای تعیین نمایش تنک سیگنال از فرم زیر استفاده می‌شود که در آن شرط تنکی با حداقل‌سازی نرم l_1 تأمین می‌شود [۲]:

$$\min \|\mathbf{s}\|_1 \quad s.t. \quad \mathbf{y} = \mathbf{Y}\mathbf{s} \quad (2)$$

از آنجا که نرم یک محدب است می‌توان کمینه‌سازی آن را با استفاده از روش‌های بهینه‌سازی محدب انجام داد؛ همچنین می‌توان آن را به صورت مسأله خطی بیان کرد و با استفاده از روش‌های برنامه‌نویسی خطی مثل الگوریتم جستجوی پایه (BP)^{۳۵} آن را حل کرد [۲۱]. از مهم‌ترین معایب این روش پیچیدگی زیاد محاسباتی و در نتیجه زمان‌بر بودن اجرای آن است؛ بنابراین اغلب با روش‌های دیگری جایگزین می‌شود.

یکی دیگر از روش‌های رایج در به دست آوردن نمایش تنک، مناسب‌ترین ارتباط یا MP^{۳۶} نام دارد [۲۲]. این روش به صورت تکراری و حریصانه عمل می‌کند، به این صورت که در هر گام تنها یکی از ضرایب تنک را مشخص می‌کند. اگر \mathbf{s} برداری \mathbf{k} -تنک باشد، \mathbf{y} می‌تواند به صورت یک ترکیب خطی از \mathbf{k} اتم دیکشنری نوشته شود. در این روش ابتدا ستون‌های مورد استفاده از ماتریس دیکشنری در ترکیب خطی آشکار شده؛ سپس ضرایب این ستون‌ها - که همان مقادیر غیرصفر بردار \mathbf{s} هستند- با حل مسأله حداقل مربعات محاسبه

^{۲۳}Over-Complete^{۲۷}Orthogonal Matching Pursuit^{۲۴}Non-Polynomial Time^{۲۸}Primal Augmented Lagrangian Methods^{۲۵}Basis Pursuit^{۲۶}Matching Pursuit^{۲۹}Dual Augmented Lagrangian Methods

۴-۲-۲- خوشه‌بندی زیرفضای تنک (SSC)

فرض کنید $\{y_i \in \mathbf{R}^D\}_{i=1}^N$ مجموعه داده‌های موجود روی n زیرفضای مستقل خطی $\{S_i\}_{i=1}^n$ با بعد $\{d_i \ll D\}_{i=1}^n$ باشند و $Y_i \in \mathbf{R}^{D \times N_i}$ مجموعه N_i داده متعلق به زیرفضای i ام باشد؛ از آنجا که نمی‌دانیم کدام داده‌ها به کدام زیرفضا تعلق دارد، ماتریس داده را بصورت:

$$\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N] = [\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_n] \in \mathbf{R}^{D \times N} \quad (7)$$

تشکیل می‌دهیم که $N = \sum_{i=1}^n N_i$ تعداد کل داده‌هاست. اگر y_i داده‌ای از زیرفضای S_i باشد، می‌تواند به صورت ترکیبی خطی از d_i نقطه از این زیرفضا نوشته شود. به عبارتی y_i نمایشی d_i -تنک خواهد داشت که به صورت پاسخ تنک معادله (۲) قابل بازیابی است و در آن برای $j \neq i$ ، $s_j \neq 0$ و $s_j = 0$ خواهد بود.

حال برای تعیین خوشه‌بندی داده‌ها به کمک نمایش‌های تنک بدست آمده، از روش خوشه‌بندی طیفی استفاده می‌شود. در این روش با استفاده از اطلاعات محلی حول هر داده، شباهت‌هایی بین جفت داده‌ها تعریف می‌شود. سپس عمل خوشه‌بندی با استفاده از ماتریس شباهت^{۳۰} داده‌ها و تقسیم داده‌ها به چندین گروه به گونه‌ای که داده‌های هر گروه با هم مشابه و داده‌های گروه‌های مختلف با یکدیگر متفاوت باشند، انجام می‌شود [۱۰].

اگر $Y_i \in \mathbf{R}^{D \times N_i - 1}$ ماتریسی باشد که از حذف ستون i ام \mathbf{Y} به دست آید، y_i نمایشی تنک بر آن خواهد داشت که می‌تواند از رابطه زیر محاسبه شود:

$$\min \|c_i\|_1 \quad s.t. \quad y_i = Y_i c_i \quad (8)$$

پاسخ بهینه $c_i \in \mathbf{R}^{N_i - 1}$ برداری است که مؤلفه‌های غیرصفر آن متناظر با ستون‌هایی از Y_i هستند که در همان زیرفضای y_i قرار دارند. پس از حل معادله فوق برای تمام نقاط $i = 1, \dots, N$ ماتریس ضرایب $\mathbf{C} = [\hat{c}_1, \hat{c}_2, \dots, \hat{c}_N]$ به دست می‌آید که در آن \hat{c}_i با صفر قرار دادن در سطر i ام c_i بدست آمده است. با استفاده از این ماتریس، گراف $G = (V, E)$ به گونه‌ای تعریف می‌شود که گوشه‌های V در آن متناظر با N نقطه داده باشند. یال $(v_i, v_j) \in E$ در صورتی وجود خواهد

می‌شود [۱۲].

از آنجا که ابتدا برچسب داده آزمون برای سیستم مجهول است، ماتریس \mathbf{Y} به صورت ترکیبی از داده‌های آموزشی از تمام c دسته در نظر گرفته می‌شود:

$$\mathbf{Y} = [\mathbf{Y}_1, \dots, \mathbf{Y}_c] = [\mathbf{y}_{1,1}, \mathbf{y}_{1,2}, \dots, \mathbf{y}_{c,N_c}] \quad (4)$$

در این رابطه هر کدام از Y_i ها زیر ماتریس‌هایی هستند که داده‌های دسته i ام را در خود جای داده‌اند. حال می‌توان با استفاده از رابطه (۲) نمایش تنک داده ورودی y را برحسب ستون‌های ماتریس دیکشنری \mathbf{Y} بدست آورد.

در SRC از میزان شباهت داده ورودی به داده‌های موجود در هر یک از دسته‌ها برای طبقه‌بندی استفاده می‌شود. بدین منظور تابع $\delta_i: \mathbf{R}^N \rightarrow \mathbf{R}^N$ بصورت زیر تعریف می‌شود و با اعمال آن بر بردار s ، N_i مؤلفه بردار s متناظر با داده‌های دسته i ام در دیکشنری حفظ شده؛ سایر مؤلفه‌ها صفر می‌شوند.

$$\delta_i(s) = [0, \dots, 0, s_{i,1}, \dots, s_{i,N_i}, 0, \dots, 0]^T \in \mathbf{R}^N \quad (5)$$

بنابراین $\hat{y}_i = \mathbf{Y} \delta_i(s)$ داده y را تنها به صورت ترکیب خطی از داده‌های دسته i ام بیان می‌کند. با اعمال این تابع بر نمایش تنک داده y و بازسازی آن با استفاده از داده‌های متعلق به هر یک از دسته‌ها در هر مرحله، در نهایت داده ورودی به دسته‌ای نسبت داده می‌شود که میزان باقی‌مانده زیر را حداقل کند:

$$\min_i r_i(y) = \|y - \mathbf{Y} \delta_i(\hat{s})\|_2 \quad (6)$$

الگوریتم زیر به طور خلاصه این روش طبقه‌بندی را بیان می‌کند:

ورودی: ماتریس دیکشنری برای c طبقه مختلف $\mathbf{Y} = [\mathbf{Y}_1, \dots, \mathbf{Y}_c] = [\mathbf{y}_{1,1}, \mathbf{y}_{1,2}, \dots, \mathbf{y}_{c,N_c}]$ داده ورودی y (الف) بهنجار کردن ستون‌های ماتریس دیکشنری \mathbf{Y} تا نرم l_2 واحد داشته باشند.

(ب) حل مسأله حداقل‌سازی نرم l_1 زیر:

$$\min \|s\|_1 \quad s.t. \quad y = \mathbf{Y}s$$

(ج) محاسبه باقیمانده $r_i(y) = \|y - \mathbf{Y} \delta_i(\hat{s})\|_2$ برای تمام طبقه‌ها $i = 1, \dots, c$

خروجی: برچسب داده y : $\text{identity}(y) = \arg \min_i r_i(y)$

³⁰Similarity Matrix

داده‌های آموزشی، به جای بازسازی سیگنال روی تمام دسته‌ها، این عمل تنها روی شبیه‌ترین داده‌های آموزشی به داده ورودی انجام شود که این داده‌ها با انجام خوشه‌بندی مشخص می‌شوند.

در این روش، داده‌های آموزشی در مرحله آموزش با استفاده از الگوریتم SSC خوشه‌بندی می‌شوند. بدین منظور ابتدا با استفاده از رابطه (۸) نمایش تنک هر یک از داده‌های آموزشی روی سایر داده‌های آموزشی به دست می‌آید. پس از تعیین ماتریس مجاورت گراف به صورت رابطه (۹) و ماتریس لاپلاسی، به صورتی که در بخش ۲-۲ توضیح داده شد خوشه‌ها مشخص می‌شوند.

در مرحله آزمایش با ورود هر داده آزمون به سیستم، ابتدا نمایش تنک آن بر داده‌های آموزشی با استفاده از رابطه (۲) به دست می‌آید. سپس در مرحله بازسازی، در هر گام ضرایب نمایش تنک متناظر با یک خوشه نگه داشته شده؛ مابقی ضرایب نمایش تنک برابر با صفر قرار داده می‌شوند. با یافتن حداقل خطای بازسازی، خوشه‌ای که داده به آن تعلق دارد شناسایی می‌شود.

در مرحله بعد برای تعیین برچسب داده مورد نظر، با استفاده مجدد از SRC، نمایش تنک داده تنها بر حسب داده‌های خوشه مربوط به آن که در مرحله قبل تعیین شده است، پیدا می‌شود. بنابراین تصمیم‌گیری نهایی در محدوده‌ای کوچکتر و تنها روی داده‌های یک خوشه انجام می‌گیرد. در شکل (۱) مراحل انجام طبقه‌بندی در روش پیشنهادی به طور خلاصه نشان داده شده است.

۶- نتایج

در این بخش ابتدا پایگاه داده استفاده شده بررسی شده؛ سپس نتایج شبیه‌سازی و تحلیل آن‌ها ارائه می‌شود.

داشت که داده y_j درنمایش تنک y_i حاضر باشد. از آنجا که اگر y_i بتواند بصورت ترکیب خطی از تعدادی نقاط یک زیرفضا شامل y_j نوشته شود، آنگاه y_j نیز می‌تواند بصورت ترکیب خطی از نقاط همان زیرفضا شامل y_i نوشته شود؛ لذا ماتریس مجاورت^{۳۱} گراف فوق بصورت $\tilde{C} = |C| + |C|^T$ خواهد بود [۳].

پس از تشکیل گراف G انتظار می‌رود تمام گوشه‌هایی که مربوط به داده‌های یک زیرفضا هستند مجموعه‌ای متصل به یکدیگر را در گراف تشکیل دهند، در حالی که گوشه‌های نمایش‌دهنده داده‌های متعلق به زیرفضاهای متفاوت هیچ یال مشترکی نداشته باشند. لذا \tilde{C} نمایشی قطری بلوکی به شکل زیر پیدا می‌کند:

$$\tilde{C} = \begin{bmatrix} \tilde{C}_1 & 0 & \dots & 0 \\ 0 & \tilde{C}_2 & \dots & 0 \\ & & \ddots & \\ 0 & 0 & \dots & \tilde{C}_n \end{bmatrix} \quad (9)$$

سپس ماتریس لاپلاسین گراف شباهت به صورت $L = D - \tilde{C}$ تشکیل می‌شود که $D \in \mathbf{R}^{N \times N}$ ماتریسی قطری است که بصورت $D_{ii} = \sum_j \tilde{C}_{ij}$ محاسبه می‌شود.

تعداد مقادیر ویژه صفر ماتریس لاپلاسین متناظر با گراف G برابر با تعداد مؤلفه‌های متصل گراف است. در حالت ایده‌ال، داده‌های متعلق به n زیرفضای مستقل خطی دارای n مؤلفه متصل در گراف شباهت خواهند بود. بنابراین، وقتی تعداد زیرفضاها معلوم نباشد می‌توان آن را به صورت تعداد مقادیر ویژه صفر ماتریس L تخمین زد. در نهایت با اعمال الگوریتم K-means به n بردار ویژه متناظر با n مقدار ویژه کوچکتر ماتریس لاپلاسین، خوشه‌بندی داده‌ها بدست می‌آید [۴].

۵- روش پیشنهادی

در ادامه با ترکیب دو روش خوشه‌بندی و طبقه‌بندی مبتنی بر نمایش تنک، روش طبقه‌بندی سلسله مراتبی^{۳۲} پیشنهاد شده است. ایده اصلی این روش، خوشه‌بندی داده‌ها و سپس انجام طبقه‌بندی در هر یک از خوشه‌های بدست آمده است. بدین صورت که پس از بدست آوردن نمایش تنک داده ورودی بر

³¹Adjacency Matrix

³²Hierarchical Classification

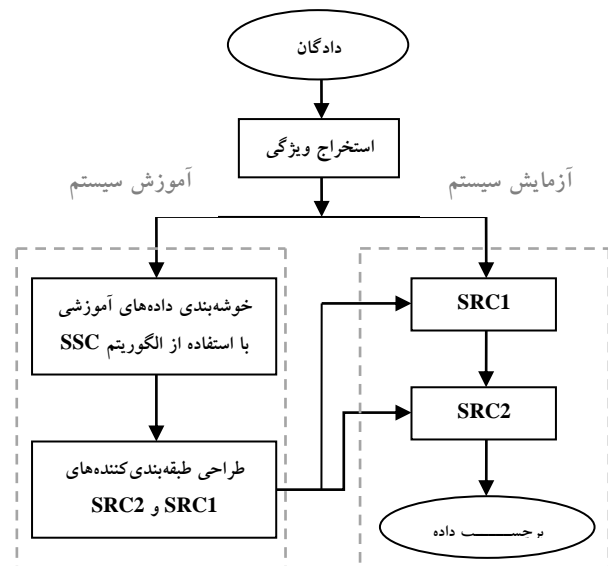
کمترین خطای به دست آمده در این روش‌ها حدود ۲۴ درصد است [۲۰]. در ادامه نتایج حاصل از اعمال روش پیشنهادی به این پایگاه داده ارائه شده است.

جدول (۱) - مشخصات دادگان مورد آزمایش

نوع سرطان	تعداد نمونه‌های آموزشی	تعداد نمونه‌های آزمون
سینه	۸	۴
پروستات	۸	۶
ریه	۸	۴
کولورکتال ^{۳۹}	۸	۴
لنفوم	۱۶	۶
مثانه	۸	۳
ملانوم ^{۴۰}	۸	۲
رحم	۸	۲
خون	۲۴	۶
کلیه	۸	۳
لوزالمعده	۸	۳
تخمدان	۸	۴
مزوتلیوما ^{۴۱}	۸	۳
سیستم عصبی مرکزی	۱۶	۴

همان‌طور که در بخش ۲ نیز ذکر شد، روش‌های متعددی برای محاسبه نمایش تنک ارائه شده است که هر یک مزایا و معایبی دارند. در جدول (۲) نتایج حاصل از استفاده از تعدادی از این روش‌ها برای محاسبه نمایش تنک در الگوریتم SRC آورده شده است. چنانچه ملاحظه می‌شود دو روش OMP و BP در این میان بهترین عملکرد را دارند و کمترین میزان خطا را در طبقه‌بندی این دادگان متحمل می‌شوند؛ اما با توجه به زمان‌بر بودن اجرای الگوریتم BP، در این کاربرد از الگوریتم OMP استفاده کرده‌ایم و نتایجی که در ادامه اعلام می‌شوند با استفاده از این الگوریتم به دست آمده‌اند.

در مرجع [۲۵] با ترکیب دو روش KNN و SRC، طبقه‌بندی انجام شده است. روش کار به این صورت است که ابتدا با اعمال طبقه‌بندی‌کننده K، KNN، نزدیک‌ترین همسایه به داده آزمون از میان داده‌های آموزشی تعیین می‌شود. سپس



شکل (۱) - روندنمای روش پیشنهادی

۷-۴-۱- دادگان مورد استفاده

برای ارزیابی عملکرد روش طبقه‌بندی معرفی شده از پایگاه داده 14-Tumors استفاده شده است [۱۶]. این مجموعه شامل ۱۹۸ نمونه حاوی ۱۶۰۶۳ ژن (ویژگی) از ۱۴ نوع سرطان است که از این میان ۱۴۴ نمونه برای آموزش سیستم و ۵۴ نمونه برای آزمایش بکار رفته‌اند. مقدار داده شده برای هر ژن در هر نمونه، میزان فعال بودن ژن در آن نمونه یا به عبارتی بیان ژن را نشان می‌دهد. مشخصات کلی این دادگان در جدول (۱) آورده شده است. همان‌گونه که از این جدول نیز ملاحظه می‌شود، تعداد نمونه در هر دسته از این داده‌ها در مقایسه با ابعاد آن‌ها بسیار کم است.

۸-۴-۲- نتایج شبیه‌سازی

از آنجا که پایگاه داده مورد استفاده در این پژوهش، از جمله پایگاه داده‌های چالش‌برانگیز در امر طبقه‌بندی است، اعمال روش‌های مختلف طبقه‌بندی بر آن اغلب نتایج ضعیفی از خود نشان داده‌اند. تا کنون از تعمیم‌های مختلفی از روش SVM شامل OVR^{۳۳}، OVO^{۳۴}، DAG^{۳۵}، WW^{۳۵} و CS^{۳۶} برای طبقه‌بندی این داده‌ها استفاده شده است. در تمام این روش‌ها از کرنل‌های چندجمله‌ای و RBF^{۳۷} استفاده شده است و انتخاب پارامترهای این طبقه‌بندی‌کننده‌ها با استفاده از روش ارزیابی متقابل^{۳۸} به صورت دقیق انجام شده است. با این حال

³³One-Versus-Rest (All)

³⁶all-at-once method by Crammer and Singer

³⁹Colorectal

³⁴Directed Acyclic Graph

³⁷Radial Basis Function

⁴⁰Melanoma

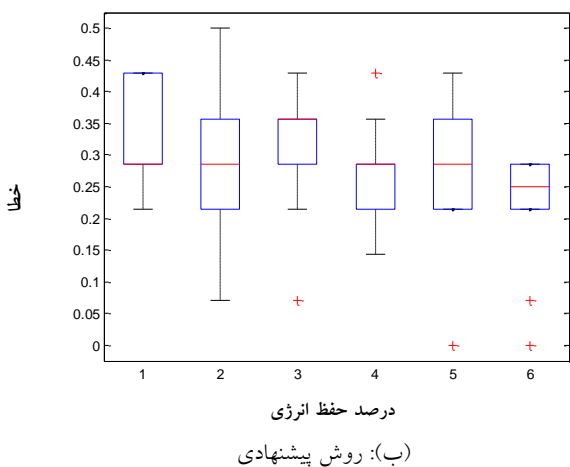
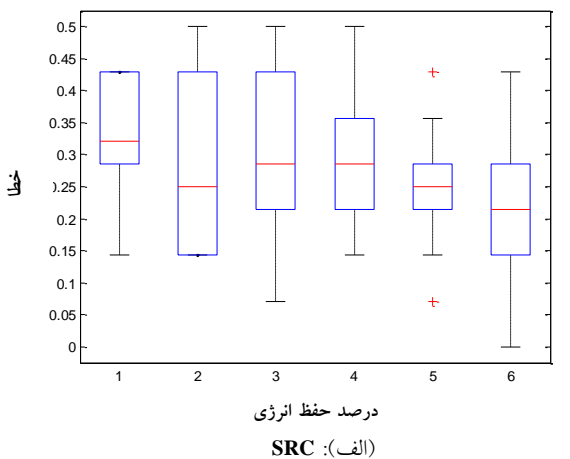
³⁵all-at-once method by Weston and Watkins

³⁸Fold Cross Validation

⁴¹Mesothelioma

است. همچنین ملاحظه می‌شود کاهش بعد داده‌ها با استفاده از PCA منجر به افزایش خطای طبقه‌بندی می‌شود.

هم در روش SRC و هم در روش پیشنهادی (SSC+SRC)، نیز به منظور کاهش بعد داده‌ها از روش استخراج ویژگی PCA استفاده شده است. به منظور تعیین بعد بهینه داده‌ها، از روش ارزیابی متقابل ۱۰-قسمتی برای داده‌های آموزشی استفاده شده است. به این ترتیب که ابتدا داده‌های آموزشی به ۱۰ قسمت مساوی تقسیم شده‌اند؛ سپس در هر مرحله یکی از قسمت‌ها به عنوان داده آزمون و ۹ قسمت باقی‌مانده به عنوان داده آموزشی برای سیستم در نظر گرفته شده‌اند. سپس عمل طبقه‌بندی به ازای شش مقدار مختلف بعد داده ورودی که به ترتیب ۹۶، ۹۷، ۹۸، ۹۹، ۹۹/۵ و ۹۹/۹ درصد انرژی داده‌های ورودی را حفظ کنند، در هر مرحله تکرار شده است. نتایج حاصل از این آزمایش‌ها در شکل (۴) نشان داده شده است.

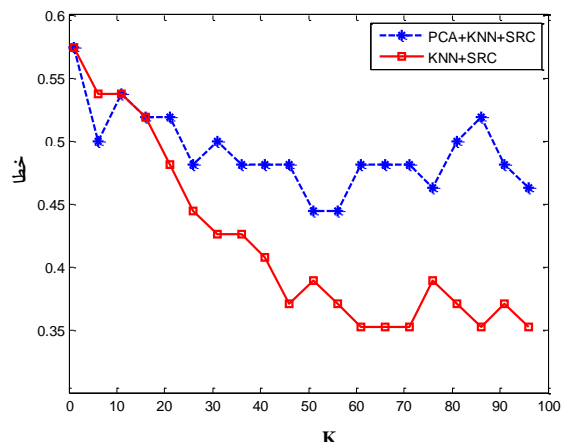


شکل (۴) - میانگین و انحراف معیار خطاهای بدست آمده در روش ارزیابی متقابل برای SRC (الف) و روش پیشنهادی (ب)

با استفاده از SRC، نمایش تنک داده آزمون تنها روی این K داده آموزشی محاسبه شده؛ طبقه‌بندی نهایی انجام می‌شود. از آنجا که ایده اصلی این روش مشابه روش پیشنهادی ما است، لذا نتایج حاصل از اعمال این روش بر دادگان مورد بررسی را نیز به منظور مقایسه با نتایج روش معرفی شده بررسی کرده‌ایم. نتایج حاصل از این بررسی در شکل (۳) ارائه شده است.

جدول (۲) - درصد خطای طبقه‌بندی با روش SRC و با استفاده از روش‌های مختلف محاسبه نمایش تنک

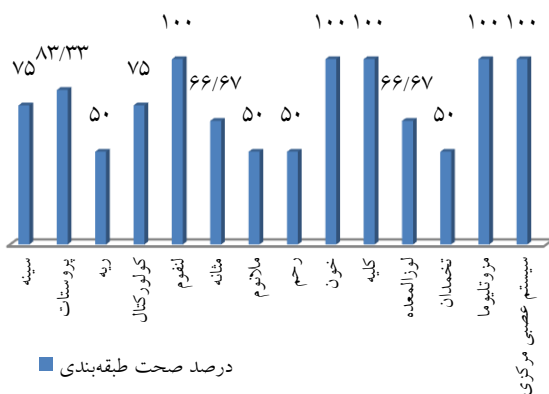
الگوریتم محاسبه پاسخ تنک	BP	OMP	Homotopy	DALM	PALM
درصد خطای طبقه‌بندی SRC با	۲۴/۰۷	۲۴/۰۷	۳۳/۳۳	۳۳/۳۳	۲۵/۹۳



شکل (۳) - خطای طبقه‌بندی در روش KNN+SRC

در این شکل مقدار خطای حاصل از طبقه‌بندی با روش KNN+SRC به ازای مقادیر مختلف K و در دو حالت بدون استخراج ویژگی و استفاده از روش استخراج ویژگی PCA آورده شده است. همان‌طور که در این شکل ملاحظه می‌شود، کمترین مقادیر خطا در حالت استفاده از تمام ویژگی‌ها و به ازای مقادیر بزرگ K بدست آمده که در حدود ۳۵ درصد

تعداد خوشه‌ها	خطای طبقه‌بندی (%)
۲	۲۴/۰۷
۳	۲۲/۲۲
۴	۲۲/۲۲
۵	۲۰/۳۷
۶	۲۰/۳۷
۷	۲۲/۲۲
۸	۲۴/۰۷
۹	۲۲/۲۲
۱۰	۲۲/۲۲
۱۱	۲۰/۳۷
۱۲	۲۰/۳۷
۱۳	۲۷/۷۸



شکل (۵) - صحت طبقه‌بندی روش پیشنهادی به تفکیک دسته‌ها

در جدول (۴) نتایج حاصل از اعمال طبقه‌بندی‌کننده‌های مختلف به پایگاه داده مورد نظر آورده شده‌اند. نتایج تعمیم‌های مختلف روش SVM و همچنین روش SR از مرجع [۲۰] آورده شده‌اند. در هر یک از این روش‌ها، بهترین نتایج اعلام شده، گزارش شده است.

با توجه به نتایج ذکر شده می‌توان دید که روش طبقه‌بندی SRC خود به تنهایی در مقایسه با روش‌های پیشین از عملکرد بهتری برخوردار است. همچنین عدم نیاز به انتخاب مدل و یا پارامترهای اولیه در مرحله آموزش را می‌توان از جمله مزایای دیگر این طبقه‌بندی‌کننده برشمرد.

در این شکل‌ها مقدار میانه و چارک‌های اول و سوم برای ده مقدار خطای بدست آمده در هر بعد، نشان داده شده است. همان‌طور که در شکل‌های بالا دیده می‌شود در هر دو روش، کمترین میزان میانه در حالت ششم و به ازای حفظ ۹۹/۹ درصد انرژی بدست آمده است. در بخش (ب) شکل (۴) حداقل تغییرات خطا نیز در این بعد بدست آمده است، ولی در قسمت (الف) این شکل در حالت پنجم کمترین تغییرات خطا رخ داده است. در مقایسه بین دو حالت پنجم و ششم برای روش SRC، از آنجایی که مقدار چارک سوم در این دو حالت برابر و میزان میانه در حالت ششم کمتر است، بعد معادل با حالت ششم را انتخاب کرده‌ایم. بنابراین بعد بهینه بمنظور اعمال الگوریتم PCA در مرحله آزمایش سیستم برای هر دو روش SRC و SSC+SRC، در ازای حفظ ۹۹/۹ درصد انرژی - که برابر با ۱۳۰ است - در نظر گرفته شده است.

نتایج طبقه‌بندی با استفاده از روش پیشنهادی به ازای داده ورودی با بعد ۱۳۰ و تعداد خوشه‌های متفاوت در جدول (۳) خلاصه شده است. هر چند در روش خوشه‌بندی SSC نیازی به دانستن تعداد خوشه‌ها نیست، در عمل به علت وجود نویز ممکن است بین داده‌های خوشه‌های مختلف نیز یال‌هایی وجود داشته باشد که تعیین تعداد خوشه‌ها را با مشکل مواجه کند؛ لذا آزمایش‌ها به ازای مقادیر مختلف تعداد زیرفضاها تکرار شده‌اند و نتایج اعلام شده است. بر اساس نتایج بدست آمده در این بررسی، در ادامه آزمایش‌ها تعداد خوشه‌ها برابر با پنج در نظر گرفته شده است.

نتایج حاصل از روش طبقه‌بندی ارائه شده به تفکیک هر یک از دسته‌ها نیز در شکل (۵) آورده شده است. همان‌طور که ملاحظه می‌شود پنج نوع سرطان لنفوم، خون، کلیه، مزوتلیوما و سیستم عصبی مرکزی به صورت کاملاً صحیح و با صحت ۱۰۰ درصد طبقه‌بندی شده‌اند و کمترین میزان صحت طبقه‌بندی مربوط به سرطان‌های ریه، ملانوم، رحم و تخمدان و به میزان ۵۰ درصد است.

جدول (۳) - خطای طبقه‌بندی روش SRC+SSC به ازای تعداد

خوشه‌های مختلف

از آنجایی که یکی از روش‌های متداول و مفید در دسته‌بندی دادگان با ابعاد زیاد و از جمله دادگان مورد بررسی، خوشه‌بندی است، لذا روش معرفی شده - که در مرحله اول خوشه مربوط به داده آزمون را تعیین کرده؛ سپس در این خوشه عمل طبقه‌بندی را انجام می‌دهد- توانسته به عملکرد قابل قبولی دست یابد. بدین ترتیب، در این روش تنها نمونه‌هایی از مجموعه آموزشی وارد فرایند طبقه‌بندی نهایی می‌شوند که بیشترین شباهت را به داده آزمون داشته باشند. بنابراین احتمال این که نمونه‌های نامربوط در نمایش تنک مربوط به داده آزمون نقش مؤثری را ایفا کنند، به شدت کاهش می‌یابد که این امر افزایش صحت طبقه‌بندی‌کننده پیشنهادی را در مقایسه با روش SRC به دنبال داشته است.

اغلب روش‌های طبقه‌بندی اعمال شده به این پایگاه داده، از تمام ۱۶۰۶۳ ویژگی برای طبقه‌بندی استفاده کرده‌اند، اما نتایج حاصل از اعمال دو طبقه‌بندی‌کننده SRC و SSC+SRC عملکرد خوبی را در ابعاد داده بسیار کمتر از بعد اولیه ویژگی‌ها نشان می‌دهند. در این روش‌ها ضمن حفظ تمام ویژگی‌ها، با انتقال آنها به فضای با بعد کمتر علاوه بر کاهش بار محاسباتی به نتایج بهتری دست یافته‌ایم. همچنین از آنجایی که روش ارائه شده روشی غیرپارامتریک است و نیازی به انتخاب مدل و تعیین پارامترهای اولیه ندارد، از این منظر نیز عملکرد آن در مقایسه با روش‌های مبتنی بر SVM قابل توجه است.

بدیهی است به منظور دستیابی به نتایج بهتر در تفکیک انواع سرطان‌ها، علاوه بر بهبود روش‌های پردازش سیگنال به استفاده بیشتر از اطلاعات بالینی مربوط به ژن‌ها و دانش زیست‌شناسی نیز نیاز است. امید است در آینده‌ای نزدیک با انجام مطالعات گسترده‌تر در این زمینه، به نتایج قابل قبولی بمنظور استفاده در تشخیص‌های بالینی دست یابیم.

مراجع

- [1] Donoho D., Compressed sensing; IEEE Trans. Information Theory, 2006; 52(4): 1289-1306.
- [2] Donoho D., For Most Large Underdetermined Systems of Linear Equations the Minimal l_1 -Norm

نتایج بدست آمده همچنین حاکی از عملکرد قابل قبول روش پیشنهادی است که منجر به کمترین مقدار خطای طبقه‌بندی در مقایسه با سایر روش‌ها شده است. اگرچه هر دو روش KNN+SRC و SSC+SRC با رویکردی یکسان نسبت به طبقه‌بندی این دادگان اقدام می‌کنند، اختلاف قابل توجهی در نتایج حاصل از آن‌ها مشاهده می‌شود که نشان‌دهنده قابلیت روش معرفی شده است. همچنین باید به این نکته نیز اشاره کرد که نتیجه روش KNN+SRC با استفاده از تمام ویژگی‌ها و به ازای مقادیر بزرگ K بدست آمده است که بار محاسباتی زیادی را به سیستم تحمیل کرده؛ زمان اجرای الگوریتم نیز طولانی می‌شود.

جدول (۴)- مقایسه نتایج حاصل از طبقه‌بندی‌کننده‌های مختلف

روش طبقه‌بندی	خطای طبقه‌بندی (%)
SVM (OVR)	۲۴/۷۱
SVM (OVO)	۵۳/۶۱
SVM (DAG)	۵۴/۹۰
SVM (WW)	۳۴/۱۶
SVM (CS)	۲۴/۶۲
SR	۲۵/۹۶
KNN+SRC	۳۵
SRC	۲۴/۰۷
SSC+SRC	۲۰/۳۷

۵- جمع‌بندی و نتیجه‌گیری

در این مقاله روش طبقه‌بندی سلسله مراتبی با ترکیب دو روش خوشه‌بندی زیرفضای تنک و طبقه‌بندی مبتنی بر نمایش تنک ارائه شد. نتایج آزمایش برای داده‌های 14-Tumors حاکی از این است که روش پیشنهادی در مقایسه با روش‌های طبقه‌بندی دیگر به صحت طبقه‌بندی بیشتری دست یافته است. بر اساس نتایج بدست آمده در جدول (۴)، نه تنها استفاده از روش SRC خطای طبقه‌بندی کمتری را در مقایسه با روش‌های پیشین در پی داشته است، بلکه ترکیب آن با روش SSC منجر به بهبود قابل توجهی در صحت طبقه‌بندی داده‌ها شده است.

- [15] Guyon I., Weston J., Barnhill S., Gene Selection for Cancer Classification using Support Vector Machines; *Machine Learning*, 2002; 46: 389–422.
- [16] Ramaswamy S., Tamayo P., Rifkin R., Multiclass cancer diagnosis using tumor gene expression signatures; *Proc. National Academy of Sciences of the United States of America*, 2001; 98(26): 15149–15154.
- [17] Linder R., Dew D., Sudhoff H., Theegarten D., Remberger K., Poppel S. J., Wagner M., The subsequent artificial neural network (SANN) approach might bring more classificatory power to ANN based DNA microarray analyses; *Bioinformatics*, 2004; 20(18): 3544–3552.
- [18] Zhang R., Huang G. B., Sundararajan N., Saratchandran P., Multi-category classification using an extreme learning machine for microarray gene expression cancer diagnosis; *IEEE/ACM Trans. on Computational Biology and Bioinformatics*, 2007; 4(3): 485–495.
- [19] Shabgahi A. Z., Abadeh M. S., A fuzzy classification system based on memetic algorithm for cancer disease diagnosis; in *Proc. 18th IEEE Iranian Conference of Biomedical Engineering (ICBME)*, 2011.
- [20] Hang X. Wu F., Sparse Representation for Classification of Tumors Using Gene Expression Data; *Journal of Biomedicine and Biotechnology*, 2009.
- [21] Chen S., Donoho D., Saunders M. A., Atomic decomposition by basis pursuit; *SIAM journal on scientific computing*, 1998; 20: 33–61.
- [22] Mallat S. G., Zhifeng Z., Matching pursuit with time-frequency dictionaries; *IEEE Trans. On Signal Processing*, 1993; 41: 3397–3415.
- [23] Pati Y. C., Rezaiifar R., Krishnaprasad P. S., orthogonal matching pursuit: recursive function approximation with applications to wavelet decomposition; 1993; rec. 27 Asilomar conf. Signals, syst. Comput, 41–44.
- [24] Yang A. Y., Zhou Z., Ganesh A., Sastry, S. S., Ma Y., Fast l_1 -Minimization Algorithms For Robust Face Recognition; *IEEE Transactions on Image Processing*, 2013; 99:1-1-0.
- [25] Nan Z., Jian Y., K nearest neighbor based local sparse representation classifier; *CCPR, IEEE*, 2010; 1–5.
- Solution Is Also the Sparsest Solution; *Comm. Pure and Applied Math.*, 2006; 59(6): 797–829.
- [3] Elhamifar E., Vidal R., Sparse subspace clustering: Algorithm, theory, and applications; to appear in *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2012; arXiv preprint arXiv: 1203.1005.
- [4] Elhamifar E., Vidal R., Sparse subspace clustering; *Proc. Int. Conf. on Computer Vision and Pattern Recognition*, 2009; 2790–2797.
- [5] Liu G., Lin Z., Yu Y., Robust subspace segmentation by low-rank representation; *Proc. Int. Conf. on International Conference on Machine Learning*, 2010.
- [6] Ho J., Yang M.H., Lim J., Lee K.C., Kriegman D., Clustering appearances of objects under varying illumination conditions; *Proc. Int. Conf. on Computer Vision and Pattern Recognition*, 2003.
- [7] Vidal R., Ma Y., Sastry S., Generalized principal component analysis (gpca); *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2005; 27(12): 1945–1959.
- [8] Tipping M.E., Bishop C.M., Mixtures of probabilistic principal component analyzers; *Neural Computation*, 1999; 11(2): 443–482.
- [9] Fischler M.A., Bolles R.C., Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography; *Communications of the ACM*, 1981; 24(6): 381–395.
- [10] Von Luxburg U., A tutorial on spectral clustering; *Statistics and Computing*, 2007; 17.
- [11] Saha B., Pham D., Phung D., Venkatesh S., Sparse Subspace Clustering via Group Sparse Coding; 2013.
- [12] Wright J., Yang A.Y., Ganesh A., Sastry S.S., Ma Y., Robust face recognition via sparse representation; *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2009; 31(2): 210–227.
- [13] Sahu S., Panda G., Barik R., A Hybrid Method of Feature Extraction for Tumor Classification Using Microarray Gene Expression Data; *International Journal of Computer Science & Informatics*, 2011; 1(1).
- [14] Sharma A., Paliwal K., Cancer classification by gradient LDA technique using microarray gene expression data; *Data & Knowledge Engineering*, 2008; 66: 338–347.