

## خوشه‌بندی مبتنی بر ماشین بردار پشتیبان دوقلو به منظور انتخاب ویژگی در مسئله

### دسته‌بندی داده‌های ریزآرایه

نفیسه سلیمانی\*<sup>(۱)</sup> محمدحسین معطر<sup>(۱)</sup>

(۱) گروه مهندسی کامپیوتر، واحد مشهد، دانشگاه آزاد اسلامی، مشهد، ایران.\*

(۲) گروه مهندسی کامپیوتر، واحد مشهد، دانشگاه آزاد اسلامی، مشهد، ایران.

تاریخ دریافت: ۱۳۹۷/۱/۲۷ تاریخ پذیرش: ۱۳۹۸/۹/۱۹

#### چکیده

طبقه‌بندی سرطان، به‌عنوان مسئله‌ای مهم در تشخیص و درمان سرطان به شمار می‌رود. یکی از مؤثرترین روش‌ها در طبقه‌بندی سرطان، شناسایی ژن‌هایی مرتبط و تبعیض‌آمیز برای طبقه‌بندی نمونه‌ها در آنالیز بیانی ژن می‌باشد. در روش پیشنهادی در این مقاله، با خوشه‌بندی ویژگی‌ها و اعمال انتخاب ویژگی درون خوشه‌ها، انتظار می‌رود که متمایزکننده‌ترین و مهم‌ترین ویژگی‌ها استخراج شوند. در روش پیشنهادی، به منظور کاهش ابعاد مجموعه داده، فن انتخاب ویژگی مبتنی بر اهمیت ویژگی‌ها به کار گرفته می‌شود، ویژگی‌های رتبه بالا استخراج شده و جهت خوشه‌بندی به ماشین بردار پشتیبان دوقلو برای خوشه‌بندی ارائه می‌شوند. پس از خوشه‌بندی، با به کار گرفتن فن انتخاب ویژگی مبتنی بر همبستگی، قابل‌اعتمادترین ویژگی‌ها انتخاب شده و توسط طبقه‌بند پرسپترون چندلایه، طبقه‌بندی می‌شوند. جهت ارزیابی روش پیشنهادی، از چهار مجموعه داده‌ی *DLBCL Leukemia*، *SRBCT* و *Prostate* استفاده شده است. نتایج آزمایش‌ها بیانگر بهبود عملکرد دقت طبقه‌بندی می‌باشد.

واژه‌های کلیدی: انتخاب ویژگی، خوشه‌بندی، ماشین بردار پشتیبان دوقلو جهت خوشه‌بندی (*TWSVC*)، طبقه‌بندی، پرسپترون چندلایه (*MLP*)

\* عهده‌دار مکاتبات:

نشانی: گروه مهندسی کامپیوتر، واحد مشهد، دانشگاه آزاد اسلامی، مشهد، ایران.

تلفن: ۰۹۱۵۵۰۳۰۵۴۰ پست الکترونیکی: [moattar@mshdiau.ac.ir](mailto:moattar@mshdiau.ac.ir)



مرتبط از فنوتایپ ها، به عنوان مثال سرطان در برابر نرمال می باشد [۵]. مهم ترین چالش ها در زمینه ی طبقه بندی سرطان عبارت اند از:

۱. با فرض اینکه  $M$  حالات ژن و  $N$  نمونه های بافتی باشد:  $M \gg N$ . محدوده ای ۲۰۰۰-۲۰۰۰۰ داشته درحالی که  $N$  محدوده ای ۳۰-۲۰۰ دارد.

۲. اغلب ژن ها برای طبقه بندی انواع مختلف بافت ها مرتبط نیستند.

۳. داده ها طبیعتی نویزی دارند.

جهت غلبه بر این چالش ها، رویکرد معمولاً انتخاب ژن جهت انتخاب زیرمجموعه ی کوچکی از ژن های آگاهی بخش که موجب به حداکثر رساندن دقت طبقه بندی توسط طبقه بندی می شود، مورد استفاده قرار می گیرد. این رویکرد مزایایی هم چون:

۱. شناسایی ویژگی های تفسیری بیشتر

۲. کاهش ابعاد داده به منظور کاهش هزینه ی محاسباتی

۳. کاهش نویز به منظور افزایش عملکرد طبقه بندی را دارد [۵].

Jain و همکارانش در [۶] مدل هیبریدی دو مرحله ای برای طبقه بندی سرطان ارائه دادند که ادغامی از انتخاب ویژگی مبتنی بر همبستگی (CFS) و بهینه سازی ازدحام ذرات بهبود یافته ی باینری (iBPSO) است. مسئله ی همگرایی به بهینه ی محلی در BPSO سنتی، توسط iBPSO کنترل می شود. Mahajan و همکارانش در [۷] جهت انتخاب ویژگی و کاهش ابعاد، روشی به نام فرآیند تحلیل سلسله مراتبی تطبیقی (A2HP) ارائه دادند. فرآیند تحلیل سلسله مراتبی یک روش تحلیل تصمیم مبتنی بر معیارهای چندگانه است که نتایج آن به دانش متخصص یا تصمیم گیرندگان وابسته است. ابتدا، پیش پردازش مجموعه بیانی ژن انجام می شود و سپس ژن های کاهش یافته به دست آمده، به عنوان ورودی برای

بهبود مداوم فنآوری بیانی ژن، توانایی اندازه گیری سطوح بیانی هزاران ژن را به طور موازی فراهم می کند [۱]. یکی از اهداف مهم تحقیقات ریزآرایه ی DNA توسعه ی ابزارهایی برای تشخیص سرطان با دقت بیشتر و مبتنی بر پروفایل ژنتیکی یک تومور می باشد. پیش بینی دقیق انواع مختلف تومور درمانی بهتر و کاهش سمی بودن را برای بیماران ممکن می سازد. انتخاب ژن موفق، به طبقه بندی انواع مختلف سرطان کمک کرده و منجر به درک بهتری از مشخصات ژنتیکی سرطان ها و استراتژی های درمانی می شود [۲].

روش های طبقه بندی مولکولی مبتنی بر الگوریتم های یادگیری ماشین در داده های ریزآرایه ی DNA بکار گرفته شده است تا نمایش دهنده ی ارتباط کلینیکی و آماری انواع مختلف تومورها باشند [۳]. یک مجموعه داده ی ریزآرایه ی معمولی بسیار پراکنده است، پراکندگی شدید و حجم کم نمونه، مانند گلوگامی جهت طبقه بندی دقیق و قدرتمند می باشد [۲]. اساساً انتخاب ویژگی باهدف انتخاب اطلاعات مرتبط و داده های حاوی اطلاعات مفید مورد استفاده قرار می گیرد. با کاهش داده ها و هرس داده های زائد و بی ربط، انتخاب ویژگی می تواند عملکرد و سرعت الگوریتم یادگیری ماشین را افزایش داده، هم چنین کاهش نیاز جهت ذخیره ی اطلاعات، ارتقاء فهم بهتر داده ها و ساده سازی تجسم داده ها را در پی دارد [۴]. شناسایی یک زیرمجموعه ی کوچک از ژن ها برپایه ی تشخیص، نه تنها دقت طبقه بندی را بهبود می بخشد بلکه دید قابل توجهی درباره ی طبیعت بیماری و مکانیزم های ژنتیکی مسئول در مقابل بیماری می یابد [۱].

وظیفه ی طبقه بندی سرطان با استفاده از داده های ریزآرایه، طبقه بندی بافت های نمونه در کلاس های

A2HP استفاده می‌شوند. A2HP از برگ خریدهای کمی و کیفی برای انتخاب ژن مای حاوی اطلاعات مفید استفاده می‌کند. Aziz و همکارانش در [۸] روشی ارائه داده‌اند که ترکیبی از رویکرد انتخاب/استخراج ویژگی برای طبقه‌بندی شبکه مای عصبی مصنوعی (ANNs) در داده مای ریزآرایه با ابعاد بالا است که از آنالیز مؤلفه مای مستقل (ICA) به‌عنوان فن استخراج و کلونی زنبورعسل مصنوعی (ABC) به‌عنوان فن بهینه‌سازی بهینه‌سازی استفاده می‌کند. Zhang و همکارانش در [۹] از ماشین بردار پشتیبان  $1-norm$  به همراه مجذور اتلاف ( $1-norm$  SVMSL) جهت اجرا سریع انتخاب ویژگی برای طبقه‌بندی سرطان استفاده می‌کنند.  $1-norm$  SVMSL، نوعی ماشین بردار پشتیبان  $1-norm$  است که در این مقاله ارائه شده که توانایی انتخاب ژن و طبقه‌بندی به‌طور هم‌زمان دارد. جهت بهبود تفسیرپذیری ژن مای منتخب برای دقت پیش‌بینی، روش انتخاب ژن بهبودیافته مبتنی بر بهینه‌سازی ازدحام ذرات باینری BPSO و اطلاعات پیشین، توسط Han و همکارانش در [۱۰] ارائه شده است. به‌منظور انتخاب ژن، BPSO اطلاعات حساسیت ژن به کلاس GCS را رمزگذاری می‌کند. اطلاعات حساسیت ژن به کلاس، که توسط ماشین یادگیری مفرط ELM از نمونه‌ها استخراج شده‌اند، در فرآیند انتخاب، طی چهار مرحله کدگذاری می‌شوند: آماده‌سازی ذرات، به‌روزرسانی ذرات، اصلاح حداکثر سرعت و اتخاذ عملیات جهش انطباقی.

در این مقاله برای ویژگی‌های مناسب و بااهمیت از خوشه‌بندی ویژگی‌ها به کمک ماشین بردار پشتیبان دوقلو استفاده گردید. جهت انتخاب آموزنده‌ترین ویژگی‌های هر خوشه فن انتخاب ویژگی مبتنی بر همبستگی استفاده شد. آنگاه با کنار هم قرار دادن ویژگی‌های مستخرج از خوشه‌های متعدد، مجموعه داده‌ی موردنظر نهایی ایجاد گردید. جهت طبقه‌بندی و

به دست آوردن دقت حاصل از طبقه‌بندی برای ویژگی‌های مستخرج و منتخب از خوشه‌های متعدد طبقه بند پرسپترون چندلایه استفاده شد.

مهم‌ترین نوآوری گروه مقاله عبارت است از خوشه‌بندی ویژگی‌های مرتبط با بیشترین مشابهت با یکدیگر و یافتن آموزنده‌ترین ویژگی‌ها در میان ویژگی‌هایی با خصیصه‌هایی مشابه از هر خوشه. نتایج حاصل از آزمایش‌ها بیانگر این امر است که انجام مراحل مذکور نه تنها منجر به کاهش چشمگیر ویژگی‌ها و پالایش مجموعه داده‌ها می‌شود، بلکه انتخاب آموزنده‌ترین ویژگی‌ها از میان ویژگی‌هایی با خصایص مشابه باعث بهبود دقت طبقه بند می‌شود.

ساختار مقاله به شرح زیر است: در بخش دوم به توصیف انتخاب ویژگی، انتخاب ویژگی مبتنی بر اهمیت ویژگی‌ها و انتخاب ویژگی مبتنی بر همبستگی؛ در بخش سوم، خوشه‌بندی و ماشین بردار پشتیبانی دوقلو برای خوشه‌بندی، و در بخش چهارم طبقه بند پرسپترون چندلایه که در روش پیشنهادی مورد استفاده قرار می‌گیرند، می‌پردازیم. در بخش پنجم، به بیان روش پیشنهادی پرداخته خواهد شد؛ ارزیابی روش پیشنهادی در بخش ششم و نتایج در بخش هفتم ارائه می‌شوند.

## ۲. انتخاب ویژگی

انتخاب ویژگی، زیرمجموعه‌ای از ویژگی‌های مرتبط را انتخاب می‌کند و ویژگی‌های زائد و بی‌ارتباط را از داده‌ها به‌منظور ساخت مدل مای یادگیری قدرتمند حذف می‌کند [۱۱]. مهم‌ترین هدف انتخاب ویژگی، کاهش ابعاد با حذف زوائد و انتخاب ویژگی‌های مرتبط می‌باشد. این امر موجب افزایش دقت یادگیری می‌شود که قابل‌فهم بودن نتایج را افزایش داده و زمان یادگیری را کاهش می‌دهد [۱۲]. در مفهوم طبقه‌بندی، فن مای انتخاب ویژگی در سه دسته با توجه به چگونگی ترکیب

آزمون مای انفرادی هستند که ویژگی‌های مرتبط با متغیر موردعلاقه (کلاس) را اندازه‌گیری می‌کنند. فرمول معادله‌ی مکاشفه، به صورت زیر می‌باشد:

$$Merit_s = \frac{\overline{k r f}}{\sqrt{k + k(k-1)r_{ff}}} \quad (1)$$

Merits مکاشفه‌ی "شایستگی" زیرمجموعه ویژگی S که حاوی k ویژگی است، میانگین همبستگی ویژگی کلاس و میانگین همبستگی ویژگی است. صورت کسر می‌تواند به عنوان نشانی داده شده از چگونگی پیش‌بینی یک گروه از ویژگی‌ها در نظر گرفته شوند؛ مخرج کسر میزان افزونگی که در میان آن‌ها وجود دارد را نشان می‌دهد. مکاشفه ویژگی‌های بی‌ارتباط را به عنوان این که در پیش‌گویی مای کلاس ضعیف عمل می‌کنند، دسته‌بندی می‌کند. ویژگی‌های اضافی در برابر آن‌هایی که با یک یا بیش از یک ویژگی همبستگی بالایی دارند، متمایز می‌شوند. به دلیل این که ویژگی‌ها به طور مستقل رفتار می‌کنند، CFS نمی‌تواند با قدرت اثرات متقابل ویژگی‌ها را شناسایی کند. با این حال می‌تواند ویژگی‌های مفید را تحت سطوح معتدلی از تعاملات شناسایی کند [۱۶].

### ۳. خوشه‌بندی

خوشه‌بندی به معنای انجام تقسیم‌بندی یک مجموعه داده‌ی بدون برچسب به گروه‌هایی با اشیاء مشابه می‌باشد [۱۷]. خوشه‌بندی به عنوان یک فن پردازش داده در حوزه مای مختلف بسیاری، همانند هوش مصنوعی، بیوانفورماتیک، بیولوژی، بینایی کامپیوتر، برنامه‌ریزی شهری، داده‌کاوی، فشرده‌سازی داده‌ها، تجزیه و تحلیل تصاویر و ... استفاده می‌شود. هر خوشه باید دو ویژگی مهم را دارا باشد؛ تشابه کم مابین

جستجو انتخاب ویژگی با ساخت مدل طبقه‌بندی سازمان‌دهی می‌شوند: روش‌های پالایه، روش‌های پوشش‌دهنده و روش‌های تعبیه شده [۱۳].

### ۱-۲. انتخاب ویژگی مبتنی بر اهمیت ویژگی‌ها

در فن انتخاب ویژگی مبتنی بر اهمیت ویژگی‌ها به هر ویژگی، مقادیر اهمیت اختصاص داده می‌شود [۱۴]. بر این اساس که اگر یک ویژگی مهم باشد، آنگاه احتمال قوی وجود دارد که عناصری با مقادیر مجموعه مای مکمل برای این ویژگی، متعلق به مجموعه‌ی مکمل از کلاس‌ها باشند. با توجه به این که تصمیمات کلاس برای دو مجموعه از عناصر متفاوت است، انتظار می‌رود که مقادیر اهمیت ویژگی برای این دو مجموعه از عناصر متفاوت باشد. اهمیت یک ویژگی به عنوان یک عملکرد دوطرفه از ارتباط با تصمیم کلاس محاسبه می‌شود. یک ویژگی در صورتی مهم است که ارتباط ویژگی با کلاس و ارتباط کلاس با ویژگی برای ویژگی بالا باشد.

### ۲-۲. انتخاب ویژگی مبتنی بر همبستگی

عوامل زیادی جهت موفقیت وظیفه‌ای محول شده در یادگیری ماشین تأثیرگذار هستند. اولین و مهم‌ترین عامل، بازنمایی و کیفیت نمونه داده‌ها می‌باشد [۱۵]. در قلب الگوریتم CFS، مکاشفه‌ای برای ارزیابی ارزش یا شایستگی یک زیرمجموعه از ویژگی‌ها وجود دارد. این مکاشفه میزان مفید بودن ویژگی‌های فردی برای پیش‌بینی برچسب کلاس به همراه سطح همبستگی میان آن‌ها را در نظر می‌گیرد. فرضیه‌ی مکاشفه به این صورت است که:

زیرمجموعه ویژگی‌های خوب شامل ویژگی‌هایی است که همبستگی بالایی با کلاس دارند در حالی که به همدیگر وابسته نیستند. در نظریه‌ی آزمون، همین اصل برای طراحی یک آزمون مرکب (مجموع یا میانگین آزمون مای انفرادی) برای پیش‌بینی یک متغیر خارجی موردعلاقه استفاده شده است. در این شرایط، ویژگی‌ها

$$\min_{w_i, b_i, x_i} \frac{1}{2} P X_i w_i + b_i e^2 + c e^T x_i \quad (3)$$

$$s.t. |X_i w_i + b_i e|^3 e - x_i, x_i^3 \geq 0$$

می‌باشد. به طور دقیق‌تر، مرکز صفحه‌ی  $i$  امین خوشه در TWSVC تا حد امکان به  $i$  امین خوشه‌ی  $X_i$  نزدیک و از دیگر خوشه‌های از هر دو سمت با  $i=1, \dots, k$  دور باشد.

با برچسب مای اولیه‌ی خوشه‌ی  $X$ ، TWSVC مراکز صفحات تمام خوشه‌ها را به‌روزرسانی می‌کند و برچسب مای نمونه‌ی جایگزین تا رسیدن به شرایط راضی‌کننده به-روزرسانی می‌شوند [۱۹].

#### ۴. طبقه‌بندی

طبقه‌بندی یکی از وظایف مهم شناسایی الگو می‌باشد [۲۰]. پرسپترون چندلایه (MLP) نوعی از شبکه‌ی عصبی است که داده مای ورودی را به داده مای هدف مورد انتظار نگاشت می‌کند. MLP متشکل از چندین لایه گره در یک گراف جهت‌دار است که هر لایه به‌طور کامل به لایه‌ی بعدی متصل است. به‌طور کلی (شکل ۱)، فرض می‌شود که  $[(x(n), t(n)]$  معرف  $n$  امین نمونه آموزشی است و  $x(n)=[x_1(n), x_2(n), \dots, x_d(n)]^T$  که  $(n=1, \dots, N)$  بردار ورودی با ابعاد  $d$  و  $t(n)=[t_1(n), t_2(n), \dots, t_c(n)]^T$  هدف با ابعاد  $c$  را نشان می‌دهند. آموزش MLP، مسئله‌ی بهینه‌سازی به حداقل رساندن مجموع میانگین مربعات خطا ( $E$  MLP بین هدف  $t_k(n)$  و خروجی واقعی  $y_k(n)$ ) می‌باشد.

$$E = \sum_{n=1}^N \sum_{k=1}^c (y_k(n) - t_k(n))^2 \quad (4)$$

با فرض این‌که  $g$  تابع فعال‌سازی در لایه پنهان،  $k$  ابعاد هدف،  $h$  تابع فعال‌سازی در لایه ورودی،  $A$  وزن مای

کلاسی و تشابه زیاد درون کلاسی. خوشه‌بندی یک یادگیری بدون نظارت می‌باشد. هیچ برچسب کلاس از پیش تعریف‌شده‌ای برای نقاط داده وجود ندارد. خوشه‌بندی موجب دستیابی به توزیع کلی الگوها و همبستگی میان اشیاء داده‌ها می‌شود [۱۸].

۱-۳. ماشین بردار پشتیبانی دوقلو برای خوشه‌بندی

روش خوشه‌بندی جدیدی مبتنی بر صفحه و بر پایه‌ی ماشین بردار پشتیبانی دوقلو (TWSVM) به نام بردار پشتیبانی دوقلو برای خوشه‌بندی (TWSVC) ارائه شده است. TWSVM نقطه‌ی عطفی در طبقه‌بندی مبتنی بر صفحه می‌باشد.

مزایای این روش عبارت‌اند از:

۱. با استفاده از کالبد TWSVM، TWSVC از اطلاعات درون و مابین خوشه‌ای بهره‌برداری می‌کند.
۲. متفاوت از TWSVM که یک صفحه‌ی کلاس برای حفظ نمونه‌های کلاس مای مختلف که به‌دوراز فقط یک صفحه هستند، موردنیاز است؛ در TWSVC، نیازمندی از فقط یک بخش به‌طور معقول‌تری از دو بخش ابر صفحه‌ی مرکز خوشه جایگزین شده است [۱۹].

۱-۳-۱. الگوریتم ماشین بردار پشتیبانی دوقلو برای خوشه‌بندی

برای مسئله‌ی خوشه‌بندی، TWSVC ارائه شده،  $k$  صفحه‌ی مرکز خوشه را جستجو می‌کند:

$$Center - plane_i = w_i^T x + b_i = 0, \quad i = 1, \dots, k \quad (2)$$

$k$  کلاس نمونه‌های  $X_1, X_2, \dots, X_k$  به‌عنوان مجموعه‌ی آموزش  $X$  می‌باشد. با در نظر گرفتن مسئله‌ی زیر با :

که  $c > 0$  یک پارامتر پنالتی و یک بردار ضعیف

حذف می‌شوند. در پایان این مرحله، تعداد ویژگی‌های مهم هر مجموعه داده تا حد ممکن حفظ‌شده و حجم هر مجموعه داده به‌طور چشم‌گیری کاهش می‌یابد.

۲- خوشه‌بندی ویژگی‌ها: در روش پیشنهادی با انجام عملیات خوشه‌بندی، ویژگی‌ها در گروه‌های سازمان‌یافته قرار می‌گیرند. ویژگی‌هایی درون خوشه‌های واحد قرار می‌گیرند که دارای حداکثر شباهت با یکدیگر و حداقل شباهت با دیگر خوشه‌ها هستند. جهت خوشه‌بندی مهم‌ترین ویژگی‌های استخراج‌شده از مرحله اول، از الگوریتم ماشین بردار پشتیبانی دوقلو برای خوشه‌بندی ((TWSVC استفاده می‌شود.

۳- انتخاب ویژگی درون خوشه‌ها، به‌منظور استخراج مهم‌ترین ویژگی‌ها: خوشه‌های متعدد ایجاد‌شده با برخورداری از خصایص مشابه درون هر خوشه دارای ارتباط ساختاری با یکدیگر هستند، از این رو با به‌کارگیری فن انتخاب ویژگی مبتنی بر همبستگی درون خوشه‌ها، می‌توان ویژگی‌های آموزنده که بیشترین همبستگی را دارند، را از میان ویژگی‌هایی با خصیصه‌های مشابه انتخاب نمود؛ قابل‌ذکر است که پس از به‌کارگیری فن انتخاب ویژگی مبتنی بر همبستگی درون هر خوشه، ویژگی‌هایی که مقدار همبستگی صفر دارند، به‌عنوان ویژگی زائد در نظر گرفته‌شده و حذف می‌شوند. سپس می‌توان تضمین نمود که ویژگی‌های منتخب و مستخرج، آموزنده‌ترین و قابل‌اعتمادترین ویژگی‌ها درون هر خوشه هستند، با کنار هم قرار گرفتن ویژگی‌های مستخرج، مجموعه داده‌ی قابل‌اعتماد نهایی تشکیل می‌شود. قابل‌ذکر است که با تغییر تعداد خوشه‌ها و افزایش آن، زمان لازم جهت خوشه‌بندی افزایش می‌یابد.

۴- طبقه‌بندی: مرحله خوشه‌بندی داده‌ها و انتخاب زیرمجموعه مناسب از ویژگی‌هایی گام میانی است و اهمیت این گام و صحت آن در گام طبقه‌بندی تعیین

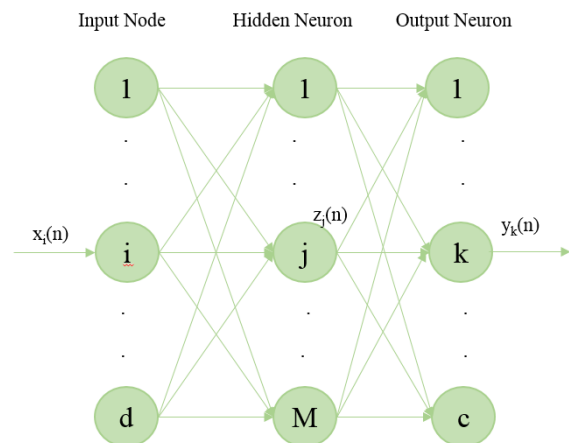
اتصال لایه‌های ورودی و پنهان و B وزن‌های اتصال لایه‌های پنهان و خروجی باشد، آنگاه داریم:

$$y_k(n) = h \sum_{j=0}^M B_{kj} z_j(n) \quad (5)$$

که  $z_j(n)$  خروجی زامین نرون در لایه پنهان را نشان می‌دهد و  $j=1,2,\dots,M$ ، با تعریف

$$z_j(n) = g \sum_{i=0}^d A_{ji} x_i(n) \quad (6)$$

MLP از دو نقص رنج می‌برد: (۱) تعیین تعداد نرون‌های پنهان دشوار است و (۲) ممکن است وزن‌آموزشی در نقاط بهینه‌ی محلی به دام بیفتند [۲۱].



شکل ۱. ساختار MLP با یک لایه پنهان [۲۱]

### ۵. روش پیشنهادی

گام‌های روش پیشنهادی به شرح ذیل می‌باشد:

- ۱- انتخاب ویژگی به‌منظور کاهش حجم ویژگی‌ها: به‌منظور کاهش حجم ویژگی‌ها، در اولین مرحله فن انتخاب ویژگی مبتنی بر اهمیت ویژگی‌ها در مجموعه داده‌ها اعمال می‌شود. تعداد بسیار زیادی از ویژگی‌های هر مجموعه داده به علت کسب رتبه پایین و مؤثر نبودن

مجموعه داده‌ی اولیه	۴	۸۳	۲۳۰۸
مرحله‌ی اول	۴	۸۳	۶۶۹
مرحله‌ی سوم	۴	۸۳	۱۶۸

جدول ۳. نتایج به دست آمده برای مجموعه داده

Leukemia

ویژگی‌ها	نمونه‌ها	کلاس‌ها	
مجموعه داده‌ی اولیه	۷۲	۲	۵۱۴۷
مرحله‌ی اول	۷۲	۲	۹۱۹
مرحله‌ی سوم	۷۲	۲	۲۵۲

جدول ۴ نتایج به دست آمده برای مجموعه داده DLBCL

ویژگی‌ها	نمونه‌ها	کلاس‌ها	
مجموعه داده‌ی اولیه	۷۷	۲	۷۰۷۰
مرحله‌ی اول	۷۷	۲	۹۰۰
مرحله‌ی سوم	۷۷	۲	۱۷۲

جدول ۵. نتایج به دست آمده برای مجموعه داده Prostate

ویژگی‌ها	نمونه‌ها	کلاس‌ها	
مجموعه داده‌ی اولیه	۱۰۲	۴	۱۲۵۳۲
مرحله‌ی اول	۱۰۲	۴	۲۲۲۱
مرحله‌ی سوم	۱۰۲	۴	۲۱۲

با توجه به نتایج ثبت شده در جدول‌های ۲، ۳، ۴، ۵، ۶

کاهش ابعاد هر مجموعه داده، کاملاً مشهود می‌باشد.

۶-۲- معیار ارزیابی عملکرد

در این مقاله، هدف ارزیابی تأثیر انتخاب ویژگی خوشه‌بندی ویژگی‌ها بر دقت طبقه‌بندی است، از این رو دقت طبقه‌بندی به عنوان مقیاس ارزیابی استفاده می‌شود.

$$accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (7)$$

در رابطه‌ی (۷) TP تعداد نمونه‌های مثبت صحیح، FP تعداد نمونه‌های مثبت کاذب، TN تعداد نمونه‌های منفی صحیح و FN تعداد نمونه‌های منفی کاذب است.

۶-۳- نتایج

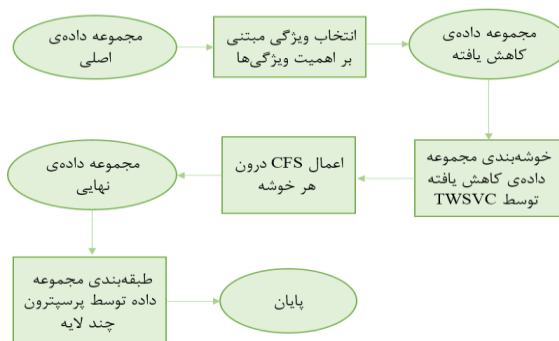
به منظور ارزیابی روش پیشنهادی، دقت طبقه‌بندی مراحل متفاوتی از اجرای روش پیشنهادی به دست آمده

می‌شود. به منظور طبقه‌بندی مجموعه داده‌ی پالایش شده و قابل اعتماد نهایی و به دست آوردن دقت و صحت الگوریتم از طبقه‌بندی پرسپترون چندلایه استفاده می‌شود.

۶. ارزیابی روش پیشنهادی

۶-۱- مجموعه داده‌ی مای ریزآرایه

جهت بررسی روش پیشنهادی، چهار مجموعه داده از وبسایت [۲۲] انتخاب شده است. خلاصه‌ای از مجموعه داده‌ها در جدول (۱) ارائه شده است.



شکل ۲. فلوچارت روش پیشنهادی

ویژگی‌ها	نمونه‌ها	کلاس‌ها	مجموعه داده
SRBCT	۸۳	۴	۲۳۰۸
Leukemia	۷۲	۲	۵۱۴۷
DLBCL	۷۷	۲	۷۰۷۰
Prostate	۱۰۲	۲	۱۲۵۳۲

جدول ۱. مشخصات مجموعه داده‌ی مورد استفاده

ویژگی‌ها نمونه‌ها کلاس‌ها مجموعه داده

نتایج به دست آمده پس از اجرای مراحل روش پیشنهادی، برای هر مجموعه داده در جدول‌های ۵، ۶، ۷، ۸ ارائه شده است.

\* دومین مرحله در روش پیشنهادی، خوشه‌بندی ویژگی‌ها می‌باشد؛ در این گام، تعداد ویژگی‌ها و نمونه‌ها تغییر نمی‌کند.

جدول ۲. نتایج به دست آمده برای مجموعه داده SRBCT

ویژگی‌ها	نمونه‌ها	کلاس‌ها	



است، نتایج در جدول ۶ ارائه شده‌اند.

ستون اول: اعمال مرحله‌ی اول از روش پیشنهادی - (Significance+MultiLayerPerceptron)

ستون دوم: اعمال مرحله‌ی اول و سوم از روش پیشنهادی

(Significance + CFS+ MultiLayerPerceptron) -

ستون سوم: روش پیشنهادی (به اختصار TWCMP) -

(Significance+TWSVC+CFS+ MultiLayerPerceptron)

با توجه به جدول ۶، می‌توان مشاهده کرد که افزودن هر یک از گام‌های روش پیشنهادی، موجب بهبود دقت طبقه‌بندی می‌شود. دقت به دست آمده برای مجموعه داده‌های SRBCT و Leukemia بدون تغییر باقی می‌ماند، اما روش پیشنهادی در مجموعه داده‌ی DLBCL، به نتیجه‌ی ۱۰۰٪ دست یافته است

جدول ۶: دقت طبقه‌بندی پیش و پس از اجرای روش پیشنهادی

Dataset	Sig+MLP	Sig+CFS+MLP	Proposed method
SRBCT	۱۰۰	۱۰۰	۱۰۰
Leukemia	۹۸.۶۱	۹۸.۶۱	۹۸.۶۱
DLBCL	۹۷.۴۰	۹۸.۷۰	۱۰۰
Prostate	۷۱.۵۶	۹۴.۱۱	۹۵.۰۹
Average	۹۱.۸۹	۹۷.۸۵	۹۸.۴۲

که نسبت به ستون اول ۲.۶٪ و نسبت به ستون دوم ۱.۳٪ بهبود داشته است. پیشرفت چشمگیر بهبود دقت طبقه‌بندی در مجموعه داده Prostate کاملاً مشهود است، روش پیشنهادی به دقت ۹۵.۰۹٪ دست یافته است، که به معنای بهبود دقت طبقه‌بندی به میزان ۲۳.۵۳٪ نسبت به ستون اول و ۰.۹۸٪ نسبت به ستون دوم می‌باشد. روش پیشنهادی به میانگین دقت ۹۸.۴۲٪ دست یافته است که مقدار بالاتری نسبت به مقادیر میانگین اعلام شده برای ستون‌های دوم و سوم را دارا می‌باشد.

از سوی دیگر، به منظور ارزیابی روش پیشنهادی با دیگر

رویکردها، دقت به دست آمده توسط روش پیشنهادی با به اختصار TWCMP، با سه روش FPRS+ISVM، CBFS+ISVM و SNR+kNN مقایسه شده است؛ نتایج حاصل در جدول ۷ و توضیح مختصری در مورد رویکردهای مذکور در ادامه ارائه شده‌اند.

اولویت‌بندی مبتنی بر مجموعه‌ی خوشن (FPRS): آنالیز اولویت، وظیفه‌ی مهمی در تصمیم‌گیری چندمعیاری است. تئوری مجموعه‌ی خوشن با جایگزینی روابط معادل با روابط غالب به حل مسئله‌ی آنالیز اولویت رسیدگی می‌کند. اولویت‌بندی مبتنی بر مجموعه‌ی خوشن، روابط اولویت‌بندی را از نمونه‌هایی که با معیارهای عددی مشخص شده‌اند، استخراج می‌کند. در واقع، روابط اولویت‌بندی در مدل مجموعه‌ی خوشن فازی گنجانده می‌شود. جهت محاسبه‌ی ارتباط میان معیار و تصمیمات، مقیاس وابستگی ویژگی مدل مجموعه‌ی فازی پاولک تعمیم داده شده است [۲۳]. یادگیری SVM استاندارد یا SVM استقرایی (ISVM)، در تلاش برای تفکیک داده‌ها در فضای ورودی با داده‌های آموزشی جدید و توسط یک ابر صفحه است؛ در مواردی که داده‌های آموزشی خطی جدایی‌ناپذیر باشند، تابع هدف یادگیری ISVM با حل مسئله‌ی بهینه‌سازی درجه دوم، ابر صفحه را به طوری که فاصله‌ی میان نزدیک‌ترین بردارها با ابر صفحه حداکثر باشد، مشخص می‌کند [۲۴]. جستجو مبتنی بر همسازی در انتخاب ویژگی (CBFS)، بر مقیاس ناهمسازی تمرکز دارد؛ به طوری که یک زیرمجموعه ویژگی با حداقل دو نمونه با مقادیر ویژگی یکسان اما بر حسب طبقات متفاوت ناسازگار محسوب می‌شود [۲۵]. نسبت سیگنال به نویز (SNR) الگوهای بیانی را با حداکثر اختلاف در میانگین حالت میان دو گروه و حداقل انحراف حالت درون هر گروه شناسایی می‌کند. نسبت سیگنال به نویز به صورت معادله‌ی زیر تعریف می‌شود:

$$SNR = \frac{m_1 - m_2}{s_1 + s_2} \quad (8)$$

به ترتیب و بیانگر میانگین و انحراف معیار از طبقه  $i$  با در نظر گرفتن ویژگی‌های متناظر می‌باشند [۲۶]. الگوریتم  $K$  همسایه نزدیک (KNN) الگوریتم ساده‌ای است که تمام موارد موجود را ذخیره کرده و موارد جدید را بر مبنای مقیاس شباهت طبقه‌بندی می‌کند.  $K$  همسایه‌ی نزدیک، در تلاش برای یافتن  $k$  نمونه با بیشترین شباهت به عنوان نزدیک‌ترین همسایگان برای نمونه‌ی داده‌شده و پیش‌بینی طبقه‌ی نمونه بر اساس اطلاعات همسایگان انتخاب شده می‌باشد. برای نمونه داده‌شده، ابتدا فاصله‌ی اقلیدسی نمونه از دیگر نقاط مجموعه داده محاسبه می‌شود. فاصله‌ها از کمترین تا بیشترین مقدار، جهت شناسایی KNN مرتب می‌شوند. با بررسی هویت کلاس در میان نقاط KNN، نظرسنجی انجام می‌شود. بر اساس کلاس اکثریت KNN، نمونه به کلاسی در مجموعه داده منتسب می‌شود. اگر کلاس منتسب یافته و کلاس واقعی نمونه تطابق داشته باشند، آزمون موفقیت‌آمیز محسوب می‌شود. بدیهی است که اکثریت آراء KNN ها، زمانی رخ می‌دهد که اکثریت متغیرهای اندازه‌گیری هم رای باشند [۲۷].

جدول ۷: مقایسه‌ی دقت طبقه‌بندی چهار روش مختلف

Dataset	TWCM P	FPRS+ ISVM	CBFS+ ISVM	SNR+kN N
SRBCT	۱۰۰	۹۸.۸۹	۹۵.۱۴	۸۱.۸۰
Leukemi a	۹۸.۶۱	۹۸.۷۵	۱۰۰	۹۴.۶۲
DLBCL	۱۰۰	۹۵.۷۱	۹۱.۴۲	۸۸.۴۵
Prostate	۹۵.۰۹	۹۳.۲۷	۹۴.۰۹	۹۲.۲۷
Average	۹۸.۴۲	۹۶.۶۵	۹۵.۱۶	۸۹.۲۸

با توجه به جدول ۷، می‌توان مشاهده کرد که در سه مجموعه داده‌ی سرطانی SRBCT، DLBCL و Prostate، روش پیشنهادی دقت بالاتری نسبت به دیگر روش‌ها داراست. در مجموعه داده‌ی سرطان خون، رویکرد

CBFS+ISVM به دقت ۱۰۰٪ دست‌یافته است که به میزان ۱.۳۹٪ عملکرد بالاتری نسبت به روش پیشنهادی دارد. به‌طور میانگین، در چهار مجموعه داده‌ی سرطانی مورد استفاده، نتایج حاصل از آزمایش‌ها، بیانگر بهبود دقت طبقه‌بندی توسط روش پیشنهادی می‌باشند؛ این بهبود به میزان ۱.۷۷٪ نسبت به روش FPRS+ISVM، 3.26٪ نسبت به روش CBFS+ISVM و ۹.۱۴٪ نسبت به روش SNR+kNN می‌باشد.

#### ۷- نتیجه‌گیری

اخیراً، فن مای زیادی به منظور انتخاب ویژگی و طبقه‌بندی داده مای ریزآرایه ارائه شده‌اند. به منظور بالا بردن دقت طبقه‌بندی، رویکرد پیشنهادی به این صورت می‌باشد: ابتدا با استفاده از روش انتخاب ویژگی مبتنی بر اهمیت ویژگی‌ها، ویژگی‌ها با اهمیت بالا تعیین می‌شوند؛ این ویژگی‌ها در مجموعه داده حفظ و ویژگی‌هایی با اهمیت کم به منظور کاهش حجم مجموعه داده حذف می‌شوند. سپس ویژگی‌های استخراج شده از مرحله‌ی اول، با استفاده از ماشین بردار پشتیبانی دوقلو برای خوشه‌بندی، خوشه‌بندی می‌شوند. در گام بعدی به منظور انتخاب مهم‌ترین ویژگی‌های درون هر خوشه، الگوریتم انتخاب ویژگی مبتنی بر همبستگی اعمال می‌شود. ویژگی‌های منتخب از تمامی خوشه‌ها در کنار یکدیگر قرار گرفته و مجموعه داده‌ی جدید، قابل اعتماد و نهایی جهت طبقه‌بندی را تشکیل می‌دهند. در مرحله‌ی آخر مجموعه داده‌ی جدید، توسط طبقه‌بند پرسپترون چندلایه، طبقه‌بندی می‌شود. نتایج به دست آمده بر روی مجموعه داده مای مورد استفاده، بیانگر کارایی روش پیشنهادی و بهبود دقت طبقه‌بندی می‌باشد.

- [1] Mohamad, Mohd Saberi, Safaai Deris, and Rosli Md Illias. "A hybrid of genetic algorithm and support vector machine for features selection and classification of gene expression microarray." *International Journal of Computational Intelligence and Applications* 5.01 (2005): 91-107.
- [2] Maulik, Ujjwal, and Debasis Chakraborty. "Fuzzy preference based feature selection and semisupervised SVM for cancer classification." *IEEE transactions on nanobioscience* 13.2 (2014): 152-160.
- [3] Yu, Hualong, and Sen Xu. "Simple rule-based ensemble classifiers for cancer DNA microarray data classification." *Computer Science and Service System (CSSS), 2011 International Conference on.* IEEE, (2011).
- [4] Xu, Qifeng, and Xuegong Zhang. "Multiclass feature selection algorithms base on R-SVM." *Signal and Information Processing (ChinaSIP), 2014 IEEE China Summit & International Conference on.* IEEE, (2014).
- [5] Mohamad, Mohd Saberi, et al. "Selecting informative genes from microarray data by using hybrid methods for cancer classification." *Artificial Life and Robotics* 13.2 (2009): 414-417.
- [6] Jain, Indu, Vinod Kumar Jain, and Renu Jain. "Correlation feature selection based improved-Binary Particle Swarm Optimization for gene selection and cancer classification." *Applied Soft Computing* 62 (2018): 203-215.
- [7] Mahajan, Shafa, and Shailendra Singh. "Informative Gene Selection Using Adaptive Analytic Hierarchy Process (A2HP)." *Future Computing and Informatics Journal* (2017).
- [8] Aziz, Rabia, et al. "Artificial neural network classification of microarray data using new hybrid gene selection method." *International Journal of Data Mining and Bioinformatics* 17.1 (2017): 42-65.
- [9] Zhang, Li, et al. "Applying 1-norm SVM with squared loss to gene selection for cancer classification." *Applied Intelligence* (2017): 1-13.
- [10] Han, Fei, et al. "A Gene Selection Method for Microarray Data Based on Binary PSO Encoding Gene-to-Class Sensitivity Information." *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)* 14.1 (2017): 85-96.
- [11] H. Liu and H. Motoda. *Computational Methods of Feature Selection*. Chapman & Hall, Boca Raton, FL, (2008).
- [12] Park, Chan Hee, and Seoung Bum Kim. "Sequential random k-nearest neighbor feature selection for high-dimensional data." *Expert Systems with Applications* 42.5 (2015): 2336-2342.
- [13] Saeys, Yvan, Iñaki Inza, and Pedro Larrañaga. "A review of feature selection techniques in bioinformatics." *bioinformatics* 23.19 (2007): 2507-2517.
- [14] Ahmad, Amir, and Lipika Dey. "A feature selection technique for classificatory analysis." *Pattern Recognition Letters* 26.1 (2005): 43-56.
- [15] Hall, Mark A. "Correlation-based feature selection of discrete and numeric class machine learning." In *Proceedings of the Seventeenth International Conference on Machine Learning*, (2000): 359-366.
- [16] Hall, Mark Andrew. "Correlation-based feature selection for machine learning." *Doctoral dissertation, University of Waikato, Dept. of Computer Science*, (1999).
- [17] P. Berkhin. *Survey of clustering data mining techniques*. In J. Kogan, C. K. Nicholas, and M. Teboulle, editors, *Grouping Multidimensional Data: Recent Advances in Clustering*. Springer, (2006): pp. 25-71.
- [18] Pratima, Depa, and Nivedita Nimmakanti. "Pattern Recognition Algorithms for Cluster Identification Problem." *Special Issue of International Journal of Computer Science & Informatics (IJCSI), Vol.-II, No.-1*, (2012).
- [19] Wang, Zhen, et al. "Twin support vector machine for clustering." *IEEE transactions on neural networks and learning systems* 26.10 (2015): 2583-2588.
- [20] Sun, Yanmin, Andrew KC Wong, and Mohamed S. Kamel. "Classification of imbalanced data: A review." *International Journal of Pattern Recognition and Artificial Intelligence* 23.04 (2009): 687-719.
- [21] Zhang, Yudong, et al. "A multilayer perceptron based smart pathological brain detection system by fractional Fourier entropy." *Journal of medical systems* 40.7 (2016): 173.
- [22] "Cancer gene expression data sets and their visualizations" [Online]. Available: <http://www.biolab.si/supp/bi-cancer/projections/>
- [23] Hu, Qinghua, Daren Yu, and Maozu Guo. "Fuzzy preference based rough sets." *Information Sciences* 180.10 (2010): 2003-2022.

- [24] Singla, Anshu, Patra, Swarnajyoti, and Bruzzone, Lorenzo. "A novel classification technique based on progressive transductive SVM learning." *Pattern Recognition Letters*, 42 (2014): 101-106.
- [25] Dash, Manoranjan, and Huan Liu. "Consistency-based search in feature selection." *Artificial Intelligence* 151.1-2 (2003): 155-176.
- [26] Begum, S., Bera, S. P., Chakraborty, D., & Sarkar, R. "Breast cancer detection using feature selection and active learning." *Proceedings of the International Conference on Advancement of Computer Communication and Electrical Technology (ACCET 2016)*. *Computer, Communication and Electrical Technology* 42 (2017): 54-59.
- [27] Mejdoub, Mahmoud, and Amar Chokri Ben. "Classification improvement of local feature vectors over the KNN algorithm." *Multimedia Tools and Applications* 64.1 (2013): 197-218.