

مِنْهُ



دانشگاه صنعتی اصفهان

دانشکده مهندسی برق و کامپیوتر

**شناسایی ژن‌های موثر در بروز بیماری با استفاده از
داده‌های ریزآرایه و آنتولوژی ژن**

پایان‌نامه کارشناسی ارشد هوش مصنوعی و رباتیک

آزاده محمدی

استاد راهنما

دکتر محمد حسین سرایی



دانشگاه صنعتی اصفهان
دانشکده مهندسی برق و کامپیوتر

پایان نامه کارشناسی ارشد رشته هوش مصنوعی و رباتیک خانم آزاده محمدی

تحت عنوان

شناسایی ژن های موثر در بروز بیماری با استفاده از داده کاوی ریز آرایه و آنتولوژی ژن

در تاریخ ۱۳۸۸/۲/۵ توسط کمیته تخصصی زیر مورد بررسی و تصویب نهایی قرار گرفت.

دکتر محمد حسین سرایی

۱- استاد راهنمای پایان نامه

دکتر منصور صالحی

۲- استاد مشاور پایان نامه

دکتر علی محمد دوست حسینی

سرپرست تحصیلات تکمیلی دانشکده

سپاس پروردگان که هر چه دارم همه از لطف اوست

بر خود لازم می‌دانم از استاد ارجمند جناب آقای دکتر سهرابی که در تمامی مراحل پایان نامه راهنما و مشوق من بودند، شکر و قدردانی نمایم، همچنین از جناب آقای دکتر صالحی به خاطر راهنموی ارزنده و حمایت‌های دلسوزانه‌شان بسیار سپاسگزارم.

از جناب آقای دکتر میرلوحی و سرکار خانم مهندس آقایانی که زحمت داوری این پایان نامه را تقبل فرمودند و نیز جناب آقای دکتر دوست حسینی رئیس محترم تحصیلات تکمیلی سپاسگزاری می‌نمایم.

همچنین در این مجال از تمامی اساتید بزرگوار دانشکده برق و کامپیوتر دانشگاه صنعتی اصفهان که افتخار نگردی آنها را داشته‌ام، شکر و قدردانی می‌نمایم.

در نهایت از خانواده عزیزم شکر می‌کنم که در همه دوران زندگی پشتیبان و همراه من بوده‌اند و همواره مدیون زحمات و لطف بی‌دریغ آنها هستم.

آزاده محمدی

اردیبهشت ۱۳۸۸

کلیه حقوق مادی مترتب بر نتایج مطالعات،
ابتکارات و نوآوریهای ناشی از تحقیق موضوع
این پایان نامه (رساله) متعلق به دانشگاه صنعتی
اصفهان است.

تقدیم بہ برترین مدیہ ہامی الہی

پدر عزیز و مادر مہربانم

فهرست مطالب

صفحه	عنوان
هشت	فهرست مطالب
۱	چکیده

فصل اول: مقدمه

۲	۱-۱ داده کاوی زیستی
۵	۲-۱ اهمیت و تاریخچه تحلیل ریزآرایه
۱۰	۳-۱ روش پیشنهادی و اهداف پایان نامه
۱۱	۴-۱ روند ارائه مطالب

فصل دوم: مروری بر مفاهیم ژنتیکی پایه

۱۲	۱-۲ مقدمه
۱۲	۲-۲ مفاهیم پایه در زیست شناسی مولکولی
۱۳	۱-۲-۲ سلول
۱۳	DNA 2-2-2
۱۶	RNA ۳-۲-۲
۱۷	4-2-2 ژن
۱۷	۵-۲-۲ پروتئین
۱۸	۶-۲-۲ بیان ژن
۲۲	3-2 فناوری ریزآرایه
۲۳	۱-۳-۲ نحوه انجام آزمایش ریزآرایه
۲۵	۲-۳-۲ چالش های داده های ریزآرایه
۲۷	۳-۳-۲ کاربردهای ریزآرایه
۲۸	۴-۲ آنتولوژی ژن

فصل سوم: تحلیل داده های ریزآرایه

۳۳	۱-۳ مقدمه
۳۴	۲-۳ مروری بر روش های پیش پردازش

1-2-3 تغییر مقیاس داده‌های بیان ژن

۲-۲-۳ مقابله با مقادیر گمشده

۳-۲-۳ نرمال‌سازی داده‌های بیان ژن

۳-۳ خوشه‌بندی داده‌های بیان ژنی

۱-۳-۳ معیار شباهت

۲-۳-۳ الگوریتم‌های خوشه‌بندی

۴-۳ رده‌بندی داده‌های بیان ژنی

۱-۴-۳ K نزدیک‌ترین همسایه

2-4-3 شبکه‌های عصبی مصنوعی

۳-۴-۳ درخت‌های تصمیم

4-۴-۳ ماشین‌های بردار پشتیبان

فصل چهارم: روش‌های انتخاب ژن

۱-۴ مقدمه

۲-۴ روش‌های انتخاب ژن

۱-۲-۴ روش‌های فیلتری

۲-۲-۴ روش‌های پوششی

3-2-4 روش‌های ادغامی

۳-۴ افزودگی در مجموعه ژن‌های انتخابی

۱-۳-۴ روش‌های مقابله با افزودگی

۴-۴ ارزیابی روش‌های انتخاب ژن

۱-۴-۴ ارزیابی بر اساس کارایی رده‌بند

۲-۴-۴ ارزیابی بر اساس عملکرد زیستی

فصل پنجم: روش پیشنهادی

1-5 مقدمه

۲-۵ تخمین مقادیر گمشده

۱-۲-۵ استفاده از الگوریتم CST در مجموعه داده‌های غیر کامل

۲-۲-۵ تلفیق اطلاعات آنتولوژی ژن با اطلاعات بیان ژن

۷۲

۳-۲-۵ تخمین مقادیر گمشده بر اساس ژن‌های مشابه

۷۳

۳-۵ انتخاب ژن

۷۳

۱-۳-۵ رتبه‌بندی ژن‌ها بر اساس معیار رتبه‌بندی

۷۵

۲-۳-۵ کاهش افزونگی با در نظر گرفتن شباهت بیانی و شباهت معنایی

۷۶

۳-۳-۵ انتخاب زیرمجموعه ژن‌ها با استفاده از روش SVMRFE

۷۸

فصل ششم: پیاده‌سازی و تحلیل نتایج

۱-۶ مقدمه

۸۰

2-6 مجموعه داده‌ها

۸۰

۱-۲-۶ مجموعه داده سرطان DLBCL

۸۱

۲-۲-۶ مجموعه داده سرطان کلون

۸۲

۳-۶ پیش‌پردازش

۸۲

۱-۳-۶ پیش‌پردازش مجموعه داده سرطان DLBCL

۸۳

۲-۳-۶ پیش‌پردازش مجموعه داده سرطان کلون

۸۷

4-6 انتخاب ژن

۸۸

۱-۴-۶ انتخاب ژن در مجموعه داده سرطان DLBCL

۸۹

۲-۴-۶ انتخاب ژن در مجموعه داده سرطان کلون

۹۶

فصل هفتم: نتیجه‌گیری و پیشنهادها

۱-۷ نتیجه‌گیری

۱۰۲

۲-۷ پیشنهادها

۱۰۵

۳-۷ دستاوردهای پژوهشی پایان‌نامه

۱۰۶

۱۰۸ مراجع

چکیده

امروزه استفاده از اطلاعات ژنتیکی افراد در تشخیص و رده‌بندی انواع بیماری‌ها از جمله سرطان‌ها، مورد توجه قرار گرفته است. یکی از بهترین و دقیق‌ترین روش‌ها در این زمینه، بررسی مقادیر بیان ژنی در افراد مختلف توسط فناوری ریزآرایه می‌باشد. یکی از مشکلات داده‌های ریزآرایه کم بودن تعداد نمونه‌ها در مقایسه با تعداد ژن‌ها است. این مسئله سبب کاهش دقت رده‌بندی و افزایش هزینه‌های محاسباتی و آزمایشگاهی می‌شود، در عین حال بسیاری از این ژن‌ها در ایجاد بیماری مورد بررسی نقشی ندارند، در نتیجه تشخیص و انتخاب ژن‌های موثر در بروز بیماری علاوه بر آنکه سبب افزایش دقت رده‌بندی و کاهش هزینه‌ها می‌شود، از نظر زیستی نیز از اهمیت ویژه‌ای برخوردار است و می‌تواند اطلاعات مفیدی درباره علل و نحوه درمان بیماری‌ها در اختیار محققین قرار دهد. تشخیص و انتخاب ژن‌های موثر در بروز بیماری، از میان هزاران ژن مورد بررسی در آزمایش ریزآرایه، انتخاب ژن نام دارد.

در این پایان‌نامه با بررسی روش‌های مختلف انتخاب ژن، تلاش شده است با بهره‌گیری از مزایای روش‌های موجود، چارچوب جدیدی برای انتخاب ژن‌های موثر در بروز بیماری ارائه شود، به نحوی که نقاط ضعف روش‌های متداول پوشش داده شوند. در روش پیشنهادی، علاوه بر داده‌های بیان ژنی از یکی دیگر از منابع معتبر موجود درباره ژن‌ها یعنی آنتولوژی ژن نیز کمک گرفته شده است. استفاده از آنتولوژی ژن در کنار مجموعه داده‌های بیان ژنی تا حدی می‌تواند محدودیت‌های ریزآرایه یعنی کم بودن تعداد نمونه‌ها و خطای احتمالی در مقادیر اندازه‌گیری شده را جبران نماید. در چارچوب ارائه شده ابتدا بخش عمده‌ای از ژن‌های غیرمرتبط با کمک روش فیلتری (فیشر) حذف می‌شوند، اما روش‌های فیلتری همبستگی موجود بین ژن‌ها را مدنظر قرار نمی‌دهند در نتیجه ژن‌های باقیمانده دارای حجم بالایی از افزونگی می‌باشند. به منظور کاهش افزونگی در ژن‌های باقیمانده، یک رویکرد حریمانه برای حذف ژن‌های مشابه پیشنهاد شده است. در این رویکرد میزان مشابهت ژن‌ها با در نظر گرفتن اطلاعات آنتولوژی ژن و داده‌های بیان ژنی و بر اساس یک معیار تلفیقی محاسبه می‌شود و سپس بر اساس این معیار، ژن‌های افزونه از مجموعه ژن‌ها حذف می‌شوند. در نهایت ژن‌های باقیمانده از این مرحله، به عنوان ژن‌های کاندید به طور دقیق‌تر توسط روش SVMRFE مورد بررسی قرار می‌گیرند تا مجموعه ژن‌های نشانگر بیماری بدست آید. روش پیشنهادی بر روی دو مجموعه داده سرطان DLBCL و سرطان کلون اعمال شده است. نتایج بدست آمده نمایانگر تاثیر مثبت روش پیشنهادی بر کارایی رده‌بندی است، به علاوه مقایسه این روش با روش‌های انتخاب ژن متداول، نشان می‌دهد که روش ارائه شده به ازای تعداد ژن‌های مساوی، از کارایی بهتری برخوردار است.

همچنین از آنجایی که بسیاری از مجموعه داده‌های ریزآرایه به دلایل مختلف از جمله وجود خراش یا گرد و غبار بر روی اسلاید، بروز خطا در حین آزمایش، اختلال در تصویر ریزآرایه و پایین بودن قدرت تفکیکی تصاویر، شامل مقادیر گمشده می‌باشند در این پایان‌نامه با استفاده از تلفیق روش خوشه‌بندی CST و آنتولوژی ژن روش نوینی برای تخمین مقادیر گمشده در مرحله پیش‌پردازش ارائه گردیده است. عملکرد روش پیشنهادی بر روی مجموعه داده سرطان DLBCL و به ازای درصد‌های مختلفی از مقادیر گمشده مورد بررسی قرار گرفته است. مقایسه نتایج حاصل از روش پیشنهادی با نتایج سایر روش‌های تخمین مقادیر گمشده، نشان می‌دهد که روش پیشنهادی می‌تواند مقادیر گمشده را با دقت بالاتری تخمین بزند.

واژه‌های کلیدی: ۱- انتخاب ژن ۲- آنتولوژی ژن ۳- بیان ژن ۴- ریزآرایه

۵- مقدار گمشده

فصل اول

مقدمه

۱ + داده‌کاوی زیستی^۱

در سال‌های اخیر، با ظهور فناوری‌های زیستی پیشرفته، حجم داده‌های زیستی به‌طور قابل ملاحظه‌ای افزایش یافته است. ذخیره‌سازی و تحلیل این حجم از داده‌ها بدون استفاده از کامپیوتر، پیچیده یا تقریباً غیرممکن است. این مسئله منجر به ایجاد یک مبحث بین رشته‌ای به نام بیوانفورماتیک^۲ شده است. بیوانفورماتیک علم نوینی است که با تلفیق علوم زیست‌شناسی، کامپیوتر و ریاضیات (به‌ویژه آمار)، تلاش می‌کند به مسائل زیستی در زمینه‌های سلولی و مولکولی پاسخ دهد [۱]. به‌طور کلی می‌توان مباحث بیوانفورماتیک را به سه دسته عمده ژنومیک^۳، پروتئومیک^۴ و ترانسکریپتومیک^۵ تقسیم نمود.

ژنومیک شامل تجزیه و تحلیل داده‌های ژنتیکی است. مجموعه اطلاعات ژنتیکی هر موجود زنده ژنوم^۶ نامیده می‌شود. اطلاعات ژنتیکی کد شده در ژنوم، خصوصیات موجودات زنده را تعیین می‌کند و از یک نسل به نسل بعدی

¹ Biological data mining

² Bioinformatics

³ Genomics

⁴ Proteomics

⁵ Transcriptomics

⁶ Genome

منتقل می‌شود. عناصر وراثتی ژن^۱ نامیده می‌شوند. به‌طور خلاصه می‌توان گفت که ژنومیک شامل توالی‌یابی^۲ و تحلیل عملکرد ژن‌ها و مقایسه ساختمان ژنی در موجودات مختلف است. تحلیل توالی DNA^۳ می‌تواند برای مقایسه اعضای مختلف از یک گونه جانداران یا به‌منظور مقایسه گونه‌های مختلف جانداران مورد استفاده قرار گیرد، همچنین تحلیل توالی‌ها می‌تواند تصویر واضح‌تری از نحوه تکامل موجودات زنده و چگونگی ارتباط آنها با یکدیگر فراهم آورد. پروژه ژنوم انسان^۴ که از سال ۱۹۹۶ تا سال ۲۰۰۳ به طول انجامید، نمونه‌ای از تحلیل توالی است، در این پروژه، کل ژنوم انسان تعیین توالی گردید و درون یک پایگاه داده قرار گرفت [۲].

رفتار سلولی و تمام فعالیت‌هایی که در سلول انجام می‌شود بر عهده پروتئین‌ها است، بنابراین برای توجیه رفتار سلولی، باید ساختار و عملکرد پروتئین‌ها شناسایی شود. پروتئین‌های هر موجود بر اساس اطلاعات ژنتیکی کد شده در ژنوم او، ساخته می‌شوند. به مجموعه پروتئین‌های یک موجود، پروتئوم^۵ گفته می‌شود. پروتئومیک شاخه‌ای از علم بیوانفورماتیک است که ساختار و عملکرد پروتئین‌ها را بررسی می‌کند و داده‌های پروتئینی را مورد تجزیه و تحلیل قرار می‌دهد. در این شاخه از علم بیوانفورماتیک، پیش‌بینی ساختار پروتئین‌ها، تعیین عملکرد پروتئین‌های ناشناخته و نحوه تعامل پروتئین‌های مختلف با یکدیگر، مورد بررسی قرار می‌گیرد [۳].

ترانسکریپتوم^۶ مجموعه کامل دستورات لازم برای ساخت پروتئین‌های مختلف در یک سلول یا موجود است. همان‌طور که بیان شد ژنوم مجموعه کامل ژن‌های یک موجود و پروتئوم مجموعه پروتئین‌های قابل تولید از ژنوم است. ترانسکریپتوم همانند یک پل ارتباطی بین ژنوم و پروتئوم است و بررسی آن مشخص می‌نماید که در هر مرحله از رشد یا تحت شرایط مختلف، کدام ژن‌ها و چگونه برای ساخت پروتئین‌های لازم، فعال شده‌اند. تحلیل این دسته از داده‌ها، ترانسکریپتومیک نامیده می‌شود. در این شاخه از علم بیوانفورماتیک تحلیل الگوهای بیان ژن^۷، چگونگی نمو و تمایز سلول‌ها و سازگاری با محیط و شرایط متغیر مورد بررسی قرار می‌گیرد. با ظهور فناوری ریزآرایه^۸، این شاخه از علم بیوانفورماتیک اهمیت زیادی پیدا کرده است. برخلاف ژنوم که تنها دربردارنده اطلاعات ایستا درباره توالی ژن‌ها است، آزمایش‌های ریزآرایه با اندازه‌گیری مقادیر بیان ژنی، اطلاعات پویایی درباره عملکرد سلولی ارائه می‌کنند [۳].

¹ Gene

² Sequence alignment

³ Deoxyribo Nucleic Acid

⁴ Human genome project

⁵ Proteome

⁶ Transcriptome

⁷ Gene expression

⁸ Microarray

به طور کلی بیوانفورماتیک سه هدف عمده را دنبال می‌کند. هدف اول عبارت است از سازمان‌دهی داده‌ها به نحوی که محققان بتوانند به اطلاعات موجود دسترسی یافته و داده‌های جدید را به این اطلاعات اضافه نمایند. هدف دوم، توسعه ابزار و منابعی است که به تحلیل داده‌ها کمک نمایند. توسعه این گونه ابزارها مستلزم مهارت در تئوری محاسبات و زیست‌شناسی است. هدف سوم، استفاده از این ابزارها برای تحلیل داده‌ها و تفسیر نتایج به صورتی است که از نظر زیستی بامعنی باشد [۱]. منظور از داده‌ها در بحث بیوانفورماتیک داده‌های زیستی مانند توالی ژنوم، داده‌های بیان ژن و توالی اسیدهای آمینه در پروتئین‌ها، می‌باشد.

به‌طور کلی بیوانفورماتیک می‌تواند به صورت کاربرد فناوری کامپیوتر در مدیریت داده‌های زیستی، تعریف شود [۳]. بیوانفورماتیک به محققین اجازه می‌دهد از پیشرفت‌های علوم کامپیوتر و آمار در تحلیل داده‌های زیستی استفاده کنند؛ با این حال هر چه حجم داده‌ها بیشتر می‌شود، استفاده از روش‌های پیچیده‌تری برای مدیریت داده‌ها، ضرورت پیدا می‌کند، بنابراین در حال حاضر هدف اصلی بیوانفورماتیک استخراج اطلاعات نهفته در این حجم انبوه از داده‌ها است. کشف این اطلاعات، می‌تواند درک روشن‌تری از پدیده‌های زیستی برای محققین فراهم آورده و منجر به ایجاد فرضیات جدیدی درباره نحوه عملکرد ژن‌ها و پروتئین‌ها، دلایل بروز بیماری‌ها و نحوه درمان آنها شود. استخراج این اطلاعات توسط روش‌های هوشمند کامپیوتری از جمله داده‌کاوی^۱ میسر می‌شود.

داده‌کاوی عبارت است از استخراج اطلاعات ناشناخته و مخفی در داده‌ها که به‌طور بالقوه مفید می‌باشند. بنابراین هدف اصلی داده‌کاوی کشف اطلاعات نهفته در داده‌ها است. پژوهش جدی روی موضوع داده‌کاوی از اوایل دهه ۹۰ شروع شد [۲]. چهار مرحله اساسی در فرایند داده‌کاوی وجود دارد که عبارتند از: جمع‌آوری داده‌ها، پیش‌پردازش^۲ داده‌ها، اعمال روش‌های داده‌کاوی و تفسیر اطلاعات. در مرحله اول، داده‌ها از منابع مختلف جمع‌آوری می‌شوند. داده‌های جمع‌آوری شده معمولاً دارای داده‌هایی با مقادیر نادرست یا نامعلوم (مقادیر گمشده^۳) می‌باشند، در نتیجه پیش از اعمال روش‌های داده‌کاوی بر روی داده‌ها، داده‌های خام باید پیش‌پردازش شوند تا مشکلات گفته شده مرتفع شود. پس از پیش‌پردازش داده‌ها، با اعمال روش‌های داده‌کاوی، الگوهای مخفی در داده‌ها شناسایی می‌شوند. در این مرحله الگوریتم‌های مختلفی می‌توانند مورد استفاده قرار گیرند. انتخاب روش مناسب، با توجه به نوع مسئله صورت می‌گیرد. در نهایت اطلاعات و الگوهای استخراج شده بر اساس معیارهای مختلف ارزیابی می‌شوند، همچنین این اطلاعات می‌تواند برای تحلیل و تفسیر، در اختیار افراد خبره در زمینه مرتبط، قرار گیرد. بررسی اطلاعات بدست آمده از داده‌کاوی توسط افراد خبره، از اهمیت ویژه‌ای برخوردار است. این کار

¹ Data mining

² Preprocessing

³ Missing value

اگرچه زمان بر است اما اطمینان می دهد اطلاعات استخراج شده دقیق و سودمند هستند، به علاوه تحلیل و تفسیر صحیح نتایج را امکان پذیر می سازد.

امروزه استفاده از روش های داده کاوی در علم بیوانفورماتیک به سرعت در حال افزایش است. حجم داده های زیستی بدست آمده از آزمایش ها و نیز وجود بانک های اطلاعاتی عظیم از این نوع داده ها، منجر به ظهور زمینه تحقیقاتی جدیدی به نام داده کاوی زیستی شده است که هدف آن پردازش و تحلیل داده های زیستی و استخراج دانش نهفته در این داده ها می باشد [۳].

۱ ۴ اهمیت و تاریخچه تحلیل ریز آرایه

یکی از مهم ترین زمینه های تحقیقاتی پزشکی، شناسایی عوامل موثر در بروز بیماری است. تعیین این عوامل، می تواند روند تشخیص و درمان بیماری ها را بهبود بخشد. بسیاری از بیماری ها منشا ژنتیکی دارند، به عنوان مثال سلول های نرمال ممکن است در اثر جهش های ژنتیکی، به سلول های سرطانی بدخیم تبدیل شوند. این تغییرات، مقادیر بیان ژن ها را تحت تاثیر قرار می دهند. بیان ژن، فرایندی است که طی آن واحدهای وراثتی یعنی ژن ها به واحدهای عملیاتی یعنی پروتئین ها تبدیل می شوند [۴].

یکی از فناوری های جدید برای اندازه گیری سطح بیان ژنی، ریز آرایه می باشد. این فناوری امکان ارزیابی تعداد بسیار زیادی از ژن ها را به طور هم زمان فراهم می کند [۵]. داده های حاصل از آزمایش های ریز آرایه می توانند اطلاعات ارزشمندی درباره نحوه عملکرد سلولی، در اختیار محققین قرار دهند. تحلیل داده های ریز آرایه یکی از مسائل مطرح در زمینه بیوانفورماتیک است که در زیرشاخه ترانسکریپتومیک قرار می گیرد. هدف از تحلیل داده های ریز آرایه بررسی تغییرات بیان ژنی در شرایط مختلف است. استخراج اطلاعات از این حجم انبوه داده ها با استفاده از روش های داده کاوی امکان پذیر می شود، از جمله استفاده از روش های خوشه بندی^۱ و رده بندی^۲ در تحلیل داده های ریز آرایه بسیار مورد توجه می باشد.

هدف از خوشه بندی داده های بیان ژن، تعیین ژن ها با الگوهای بیانی مشابه و تشخیص الگوهای مخفی درون داده ها است. خوشه بندی برای تعیین ژن های مشابه، پیش بینی عملکرد ژن های ناشناخته و کاهش پیچیدگی مورد استفاده قرار می گیرد [۶]. با توجه به اهمیت شناسایی الگوهای بیان ژنی، تا کنون روش های متفاوتی برای خوشه بندی داده های بیان ژنی مورد استفاده قرار گرفته اند که از آن جمله می توان به خوشه بندی سلسله مراتبی^۳ [۷، ۸] خوشه بندی

^۱ Clustering

^۲ Classification

^۳ Hierarchical

K میانگین^۱ [۶] و نقشه خودسازمان ده^۲ [9] اشاره نمود. علاوه بر روش های مرسوم خوشه بندی، الگوریتم هایی برای خوشه بندی داده های بیان ژنی به طور خاص ارائه شده اند که از آن جمله می توان به الگوریتم CAST^۳ [10] و CST^۴ [۱۱] اشاره نمود. این روش ها در برابر نویز^۵ بهتر عمل کرده و قادر به شناسایی خودکار تعداد خوشه ها می باشند. به طور کلی بررسی ها نشان داده اند که الگوریتم های خوشه بندی که به طور خاص برای داده های بیان ژنی طراحی شده اند، در اکثر موارد نسبت به الگوریتم های متداول خوشه بندی مانند خوشه بندی K میانگین، SOM و غیره موفق تر عمل می نمایند [۶].

یکی دیگر از کاربردهای داده کاوی در ریزآرایه، رده بندی نمونه ها بر اساس مقادیر بیان ژنی است. با داده کاوی ریزآرایه، می توان الگوهایی جهت تشخیص بیماری ها استخراج نمود. تشخیص زود هنگام بیماری باعث افزایش احتمال بهبود و کاهش هزینه های درمانی می گردد، به علاوه تشخیص بیماری با استفاده از الگوهای استخراج شده از داده های ریزآرایه دارای دقت بالایی است. بسیاری از بیماری ها علائم ظاهری و بالینی مشابهی را بروز می دهند، در حالی که ممکن است در سطح سلولی متفاوت بوده و نیاز به رویکردهای درمانی مختلفی داشته باشند. با بررسی بیماری در سطح سلولی می توان گونه های مختلف یک بیماری را تفکیک و رده بندی نمود. تا کنون روش های مختلفی برای رده بندی داده های بیان ژنی مورد استفاده قرار گرفته است که از آن جمله می توان به روش K نزدیک ترین همسایه^۶ [۱۲]، شبکه های عصبی مصنوعی^۷ [۱۳]، درخت های تصمیم^۸ [۳] و ماشین های بردار پشتیبان^۹ [۱۴، ۱۵] اشاره کرد.

همان گونه که بیان شد داده های بیان ژنی می توانند در تشخیص بیماری ها مورد استفاده قرار گیرند، اما مشکل اساسی در مورد این داده ها، محدود بودن تعداد نمونه های مورد آزمایش در مقایسه با تعداد ژن ها است. بسیاری از این ژن ها در ایجاد بیماری مورد بررسی نقشی ندارند، از این رو تشخیص ژن های موثر در بروز بیماری می تواند علت و نحوه ایجاد بیماری را مشخص نموده و پزشکان و داروسازان را در اتخاذ روش های درمانی مناسب و طراحی داروهای مفید، یاری کند. با تعیین ژن های موثر در بروز بیماری، آزمایش های تشخیصی و رده بندی بیماران، می تواند

¹ K-means

² Self Organizing Map

³ Cluster Affinity Search Technique

⁴ Correlation Search Technique

⁵ Noise

⁶ K Nearest Neighbor

⁷ Artificial Neural Network

⁸ Decision tree

⁹ Support Vector Machines

تنها با استفاده از این ژن‌ها صورت گیرد؛ این امر سبب کاهش هزینه‌های آزمایشگاهی و افزایش دقت می‌شود. تشخیص و انتخاب زیرمجموعه‌ای از ژن‌ها به‌عنوان عوامل موثر در بروز بیماری، انتخاب ژن نامیده می‌شود.

تا کنون روش‌های مختلفی برای تعیین ژن‌های موثر در بروز بیماری، ارائه شده‌اند. روش‌های انتخاب ژن به‌طور کلی در سه دسته فیلتری^۱، پوششی^۲ و ادغامی^۳ قرار می‌گیرند [۱۶]. در روش‌های فیلتری، انتخاب ژن مستقل از رده‌بندی انجام می‌شود. این روش‌ها معمولاً بر پایه روش‌های آماری استوار هستند، بدین ترتیب که قدرت هر ژن در تفکیک نمونه‌ها بر اساس یک معیار آماری محاسبه می‌شود و سپس ژن‌هایی که بر اساس معیار محاسبه شده قدرت تفکیک بهتری دارند، به‌عنوان مجموعه‌ژن‌های موثر در بروز بیماری انتخاب می‌شوند. از جمله این روش‌ها می‌توان به انتخاب ژن با استفاده از معیار SNR^۴ [۱۷]، انتخاب ژن با استفاده از معیار فیشر^۵ [۱۸] و انتخاب ژن با استفاده از معیار مجموع رتبه ویلکوکسون^۶ [۱۹] اشاره کرد. این روش‌ها ساده و سریع هستند اما همبستگی بین ژن‌ها را در نظر نمی‌گیرند و مجموعه ژن‌های انتخابی دارای حجم بالایی از افزونگی^۷ می‌باشند.

در روش‌های انتخاب ژن پوششی، فضای ویژگی^۸ (ژن‌ها) برای یافتن مجموعه ژن‌های موثر در بروز بیماری جستجو می‌شود. در این روش‌ها با استفاده از یک مکانیزم جستجو، در هر مرحله یک زیرمجموعه ژن انتخاب شده و کیفیت آن بر اساس کارایی رده‌بندی^۹، مورد ارزیابی قرار می‌گیرد. زیرمجموعه‌ای از ژن‌ها که بالاترین کارایی را در رده‌بندی ایجاد نماید به‌عنوان مجموعه ژن‌های متمایزکننده انتخاب می‌شود. به‌عنوان مثال در [۲۰] از روش انتخاب رو به جلو^{۱۰} برای انتخاب ژن‌های موثر در بیماری استفاده شده است. در [۲۱, ۲۲] نیز، از الگوریتم ژنتیک^{۱۱} برای انتخاب ژن استفاده شده است. این روش‌ها تعامل ژن‌ها را در نظر می‌گیرند اما پیچیدگی محاسباتی بالایی دارند.

دسته دیگر از روش‌های انتخاب ژن روش‌های ادغامی هستند که در آنها مرحله انتخاب ژن با مرحله رده‌بندی ادغام می‌شود. روش SVMRFE^{۱۲} یکی از موفق‌ترین روش‌های انتخاب ژن است که توسط گیون^{۱۳} و همکاران ارائه شده است [۲۳]. در این روش در هر مرحله کم‌اهمیت‌ترین ژن‌ها در تفکیک نمونه‌ها، به‌صورت بازگشتی حذف

¹ Filter

² Wrapper

³ Embedded

⁴ Signal to Noise Ratio

⁵ Fisher

⁶ Wilcoxon rank sum test

⁷ Redundancy

⁸ Feature sapce

⁹ Classifier

¹⁰ Forward selection

¹¹ Genetic algorithm

¹² Support Vector Machine Recursive Feature Elimination

¹³ Guyon

می‌شوند. تعیین اهمیت ژن‌ها در این روش بر اساس ماشین بردار پشتیبان صورت می‌گیرد. این روش تعامل بین ژن‌ها را در نظر می‌گیرد و در عین حال پیچیدگی محاسباتی آن نسبت به روش‌های پوششی کمتر است.

هدف از انتخاب ژن، حذف ویژگی‌های غیرمرتبط^۱ و افزونه^۲ می‌باشد. ویژگی‌های غیرمرتبط ویژگی‌هایی هستند که در تفکیک نمونه‌ها نقش چندانی ندارند و بدون تاثیر گذاشتن بر کارایی یادگیری می‌توان آنها را از مجموعه ژن‌ها حذف کرد. ویژگی افزونه خصیصه‌ای است که به‌تنهایی مرتبط با تفکیک رده‌ها است اما وجود یک ویژگی دیگر در مجموعه ژن‌ها سبب می‌شود که حذف این ویژگی تغییری در کارایی یادگیری ایجاد نکند. اکثر روش‌های مطرح شده در انتخاب ژن به دنبال حذف ویژگی‌های غیرمرتبط هستند اما اکثر روش‌ها بخصوص روش‌های فیلتری افزونگی ژن‌ها را در نظر نمی‌گیرند. وجود افزونگی در مجموعه ژن‌های انتخابی علاوه بر آن که سبب افزایش هزینه محاسباتی می‌شود، کارایی رده‌بندی را نیز کاهش می‌دهد، همچنین افزونگی ژن‌ها سبب انتخاب ژن‌هایی با عملکرد مشابه می‌شود. این ژن‌ها حاوی اطلاعات مشابهی هستند، در نتیجه انتخاب ژن‌های افزونه، اطلاعات بیشتری را فراهم نمی‌کند [۲۴].

تا کنون روش‌های مختلفی برای کاهش افزونگی در مجموعه ژن‌های انتخابی پیشنهاد شده است. از جمله در [۲۵]، [۲۶] استفاده از روش‌های خوشه‌بندی برای تعیین ژن‌های مشابه و کاهش افزونگی پیشنهاد شده است. در [۲۵] پیش از اعمال روش‌های رتبه‌بندی، ژن‌های مشابه با استفاده از خوشه‌بندی فازی^۳ تعیین و حذف می‌شوند و رتبه‌بندی ژن‌ها تنها بر روی ژن‌های باقیمانده انجام می‌شود، اما در [۲۶]، خوشه‌بندی سلسله‌مراتبی پس از انتخاب ژن توسط روش‌های فیلتری رتبه‌بندی اعمال شده است تا افزونگی موجود در ژن‌های انتخابی را کاهش دهد. در [۲۷] از روش مارکوف بلانکت^۴ برای تشخیص و کاهش افزونگی ژن‌ها استفاده شده است. دینگ^۵ و پنگ^۶ در [۲۸] روشی به نام MRMR^۷ برای انتخاب ژن ارائه داده‌اند که هدف آن انتخاب مجموعه ژنی است که به‌طور هم‌زمان کمترین میزان افزونگی و بیشترین میزان ارتباط با رده^۸ مورد بررسی را داشته باشد.

خطاهای آزمایشگاهی در روند انجام آزمایش ریزآرایه سبب ایجاد مقادیر نادرست یا نامعلوم (مقادیر گمشده) در داده‌های بیان ژنی می‌شوند، لذا پیش از اعمال روش‌های داده‌کاوی، لازم است داده‌ها با استفاده از روش‌های پیش‌پردازش اصلاح شوند. مراحل پیش‌پردازش وابسته به خصوصیات مجموعه داده است. در بسیاری از مجموعه

¹ Irrelevant

² Redundant

³ Fuzzy clustering

⁴ Markov Blanket

⁵ Ding

⁶ Peng

⁷ Minimum Redundancy Maximum Relevance

⁸ Class

داده‌ها، مقادیر بیان ژنی در برخی از نمونه‌ها نامعلوم (گمشده) است. اکثر روش‌های داده‌کاوی برای عملکرد صحیح نیاز به مجموعه کاملی از داده‌ها دارند، از این رو مقادیر گمشده باید به شکل صحیح مدیریت شوند، همچنین به منظور حداقل کردن تاثیرات ناشی از خطای آزمایش، بر مقادیر بیان ژنی باید از روش‌های نرمال‌سازی استفاده نمود [۵].

تا کنون برای مقابله با مقادیر گمشده روش‌های مختلفی ارائه شده‌اند. از جمله در [۷] جایگزینی مقادیر گمشده با یک مقدار ثابت مانند صفر یا میانگین مقادیر ژن در نمونه‌ها پیشنهاد شده است. این روش‌ها ساده و سریع هستند اما به علت نادیده گرفتن همبستگی بین داده‌ها نتایج دقیقی تولید نمی‌کنند. ترویانسکایا^۱ و همکاران در [۲۹] روش نسبت‌دهی K نزدیک‌ترین همسایه (KNNImpute)^۲ را برای تخمین مقادیر گمشده پیشنهاد داده‌اند. در این روش مقادیر گمشده هر ژن با توجه به مقادیر K ژنی که بیشترین میزان مشابهت را با آن دارند، تخمین زده می‌شود. در [۳۰] روشی به نام نسبت‌دهی K نزدیک‌ترین همسایه متوالی (SKNNImpute)^۳ برای تخمین مقادیر گمشده ارائه شده است. این روش مشابه با روش K نزدیک‌ترین همسایه می‌باشد با این تفاوت که انتخاب نزدیک‌ترین همسایه‌ها برای هر ژن حاوی مقدار گمشده، از بین ژن‌های فاقد گمشدگی صورت می‌پذیرد و پس از تخمین مقادیر گمشده یک ژن، آن ژن نیز می‌تواند در مراحل بعدی در تخمین مقادیر گمشده سایر ژن‌ها مورد استفاده قرار گیرد. کیم^۴ و همکاران در [۳۱] روشی به نام نسبت‌دهی حداقل مربعات محلی (LLSImpute)^۵ ارائه داده‌اند که در آن مقادیر گمشده بر اساس ترکیب خطی K نزدیک‌ترین همسایه تخمین زده می‌شوند. در [۳۲, ۳۳] از روش‌های مبتنی بر خوشه‌بندی برای تخمین مقادیر گمشده استفاده شده است. در این روش‌ها، با استفاده از خوشه‌بندی، ژن‌های مشابه تعیین شده و سپس مقادیر گمشده هر ژن با استفاده از مقادیر متناظر در ژن‌های هم‌گروه آن تخمین زده می‌شوند.

مشکل عمده روش‌های فوق وابسته بودن آنها به پارامترهای ورودی مانند تعداد همسایه‌ها یا تعداد خوشه‌ها است. همچنین تخمین مقادیر گمشده در تمامی روش‌های مذکور تنها بر اساس اطلاعات بیان ژن صورت می‌گیرد، در حالی که به علت کمی نمونه‌ها و وجود خطا در داده‌های ریزآرایه تخمین مقادیر گمشده بر پایه این داده‌ها نمی‌تواند چندان دقیق باشد. در این حالت بهره‌گیری از سایر منابع اطلاعاتی موجود، می‌تواند سبب بهبود تخمین شود.

¹ Troyanaskaya

² K Nearest Neighbor Imputation

³ Sequential KNNImpute

⁴ Kim

⁵ Local Least Squares Imputation

۱ ۳ روش پیشنهادی و اهداف پایان نامه

همان گونه که بیان شد، از داده‌های بیان ژنی می‌توان در تشخیص و رده‌بندی بسیاری از بیماری‌ها مانند سرطان استفاده کرد، اما این داده‌ها دو مشکل عمده دارند: یکی آن که داده‌های حاصل از آزمایش ریزآرایه در اثر خطاهای آزمایشگاهی ممکن است مقادیر غیرصحیح یا نامعلوم (مقادیر گمشده) داشته باشند. از این رو پیش‌پردازش داده‌های بیان ژنی به منظور رفع این مشکلات پیش از انجام رده‌بندی ضروری است. مشکل دیگر داده‌های بیان ژنی، بیشتر بودن تعداد ژن‌ها در مقایسه با تعداد نمونه‌ها است که سبب کاهش دقت رده‌بندی و افزایش هزینه‌های محاسباتی و آزمایشگاهی می‌شود. در این پایان‌نامه سعی می‌شود با توجه به نقاط ضعف و قوت روش‌های موجود، روش‌های جدیدی برای برخورد با این مسائل ارائه شود.

هر چند فناوری ریزآرایه حاوی اطلاعات بسیار ارزشمندی می‌باشد، اما به دلیل امکان بروز خطا در مراحل مختلف آزمایش و آماده‌سازی داده‌ها و نیز به علت کم بودن تعداد نمونه‌های آزمایشی، تصویر کاملی از نحوه عملکرد ژن‌ها ارائه نمی‌کند؛ از این رو بهره‌گیری از سایر منابع اطلاعاتی موجود درباره ژن‌ها، می‌تواند تا حدی محدودیت‌های ریزآرایه را جبران نماید. یکی از مهم‌ترین منابع اطلاعاتی موجود درباره ژن‌ها، آنتولوژی ژن^۱ است که به شکل ساخت‌یافته عملکرد ژن‌ها و چگونگی ارتباط آنها را توصیف می‌کند. تلفیق این اطلاعات با داده‌های حاصل از آزمایش ریزآرایه می‌تواند نتایج قابل‌اطمینان‌تری تولید نماید، زیرا در این حالت تحلیل عملکرد ژن‌ها بر اساس منابع اطلاعاتی گسترده‌تری صورت می‌گیرد.

در این پایان‌نامه سعی می‌شود با تلفیق اطلاعات آنتولوژی ژن و داده‌های بیان ژن، مقادیر گمشده در داده‌های ریزآرایه تخمین زده شوند. در روش ارائه شده، با استفاده از خوشه‌بندی CST و نیز بهره‌گیری از اطلاعات آنتولوژی ژن، ژن‌های مشابه در قالب خوشه‌ها گروه‌بندی می‌شوند، سپس مقادیر گمشده هر ژن با توجه به مقادیر ژن‌های هم‌دسته آن تخمین زده می‌شوند. با این روش علاوه بر داده‌های حاصل از آزمایش ریزآرایه می‌توان از دانش موجود درباره ژن‌ها، که به صورت آنتولوژی ژن وجود دارد نیز بهره گرفت، به علاوه این روش نیازی به تعیین پارامتر توسط کاربر ندارد و تعداد ژن‌های مورد نیاز برای تخمین مقادیر هر ژن به‌طور خودکار تعیین می‌شود.

همچنین در این پایان‌نامه یک چارچوب جدید جهت شناسایی و انتخاب ژن‌های موثر در بروز بیماری ارائه می‌شود. در این روش ابتدا با استفاده از معیار فیشرفیلتر بخشی از ژن‌ها حذف می‌شوند، سپس با بهره‌گیری از اطلاعات آنتولوژی ژن و نیز داده‌های بیان ژنی، همبستگی بین ژن‌ها بررسی می‌شود و با استفاده از رویکردی حریمانه^۲، ژن‌های افزونه از مجموعه ژن‌های باقیمانده حذف می‌شوند. در نهایت با استفاده از روش SVMRFE بهترین

¹ Gene Ontology

² Greedy

زیرمجموعه ژنی انتخاب می‌شود. این ژن‌ها به‌عنوان مجموعه ژن‌های موثر در بروز بیماری می‌توانند در رده‌بندی نمونه‌ها مورد استفاده قرار گیرند. کارایی روش پیشنهادی بر روی مجموعه داده‌های سرطان DLBCL^۱ و کلون^۲ مورد بررسی قرار گرفته است.

۱ ۴ روند ارائه مطالب

در فصل دوم، مقدمه مختصری در مورد علم ژنتیک و مفاهیم پایه زیست‌شناسی مولکولی بیان شده و سپس مفهوم بیان ژنی مطرح می‌گردد، در ادامه فناوری ریزآرایه، نحوه استخراج داده‌های بیان ژنی و کاربردهای آن شرح داده می‌شود. همچنین در این فصل آنتولوژی ژن به‌عنوان یک منبع اطلاعاتی مفید در تحلیل ژن‌ها، معرفی می‌شود.

در فصل سوم، مراحل تحلیل داده‌های ریزآرایه مورد بررسی قرار می‌گیرد. در این فصل ابتدا مروری بر روش‌های پیش‌پردازش داده‌های ریزآرایه مانند روش‌های تخمین مقادیر گمشده و نرمال‌سازی داده‌ها می‌شود، سپس انواع روش‌های مطرح در خوشه‌بندی و رده‌بندی داده‌های بیان ژنی بررسی می‌شوند.

در فصل چهارم، اهداف و دلایل انتخاب ژن مورد بررسی قرار گرفته و انواع روش‌های انتخاب ژن و مزایا و معایب هریک بررسی می‌شود، همچنین مشکلات ناشی از وجود افزونگی در مجموعه ژن‌های انتخابی، مطرح شده و تعدادی از مهم‌ترین روش‌های مقابله با افزونگی معرفی می‌شوند. در پایان این فصل چگونگی ارزیابی روش‌های انتخاب ژن مورد بررسی قرار می‌گیرد.

در فصل پنجم یک روش جدید برای تخمین مقادیر گمشده ارائه می‌شود. روش ارائه شده یک روش نسبت‌دهی مبتنی بر خوشه‌بندی است که در آن برای خوشه‌بندی ژن‌ها علاوه بر داده‌های بیان ژن از اطلاعات آنتولوژی ژن نیز کمک گرفته می‌شود. همچنین یک چارچوب جدید برای شناسایی و انتخاب ژن‌های موثر در بروز بیماری ارائه می‌شود. این چارچوب با ترکیب روش‌های فیلتری و ادغامی از مزایای هر دو روش بهره می‌گیرد، همچنین با تلفیق اطلاعات آنتولوژی ژن و داده‌های بیان ژنی، افزونگی در مجموعه ژن‌های انتخابی را کاهش می‌دهد.

نتایج پیاده‌سازی روش‌ها بر روی دو مجموعه داده سرطان DLBCL و سرطان کلون در فصل ششم ارائه می‌شود. در نهایت در فصل هفتم نتایج بدست آمده جمع‌بندی شده و پیشنهادهایی جهت ادامه کار ارائه می‌شود. لیست مقالات برگرفته شده از این پایان‌نامه در انتهای فصل هفتم آورده شده است.

¹ Diffuse Large B-Cell Lymphoma

² Colon