

دانشگاه یزد

دانشکده ریاضی

گروه علوم کامپیوتر

پایان نامه

جهت دریافت درجهی کارشناسی ارشد

علوم کامپیوتر

تحلیل خوشه‌بندی برای داده بیان ژن با استفاده از تجزیه‌ی ماتریس نامنفی

استاد راهنما:

دکتر سید ابوالفضل شاهزاده‌فاضلی

استاد مشاور:

دکتر مهدیه هاشمی‌نژاد

پژوهش و نگارش:

فاطمه السادات فاطمیون

مهر ۱۳۹۰

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

سپاس‌گزارى

سپاس خدا را که مرا به راه دانش رهنمون شد و هموست که فرصت شاگردى اساتید گرامى را بر من ارزانى داشت تا چراغى بر تاریكى جهلم باشند.

از استاد بزرگوارم، جناب آقای دکتر سید ابوالفضل شاهزاده‌فاضلى، استاد راهنمای این پژوهش، که افق‌های جدیدى بر پنجره‌ى ذهنم گشودند، سپاس‌گزارم و رهنمودهای ایشان را ارج مى‌نهم.

از خانم دکتر مهدیه هاشمى‌نژاد که زحمت مطالعه و مشاوره این پایان‌نامه را تقبل فرمودند تقدیر و تشکر مى‌کنم.

از مادرم که پشتیبان همیشگى و پدرم که بزرگ‌ترین مشوق من در راه فراگیرى دانش بوده‌اند کمال سپاس و امتنان را دارم و توفیق جبران هر قطره از دریای بی‌کران مهر این عزیزان را غنیمت مى‌شمارم.

چکیده

امروزه حجم عظیمی از مطالعات پزشکی در جهت شناسایی و درمان بیماری‌هایی است که از طریق ژن منتقل می‌شود. برای بررسی و نگهداری اطلاعات ژنتیکی، فناوری‌های مفیدی به وجود آمده است که یکی از آن‌ها، فناوری ریزآرایه می‌باشد. تجزیه و تحلیل اطلاعات به دست آمده از ریزآرایه‌ها به کمک روش‌های داده‌کاوی انجام می‌شود. یکی از این روش‌ها خوشه‌بندی است که می‌تواند در یافتن گروه‌های واقعی و نهفته در داده‌ها مؤثر باشد. همچنین با استفاده از روش‌های کاهش بعد می‌توان مجموعه داده‌هایی با حجم کوچک‌تر از مجموعه داده‌های اصلی تولید کرد و آن را به‌عنوان ورودی روش خوشه‌بندی به کار برد.

در این رساله از تجزیه‌ی ماتریس نامنفی (NMF) برای کاهش بعد داده‌های ریزآرایه استفاده می‌شود. همچنین برای مقداردهی اولیه این تجزیه‌ی روش‌های تصادفی، تحلیل مولفه اصلی (PCA) و تجزیه‌ی مقدار تکین مضاعف نامنفی (NNDSVD) به کار می‌رود. پس از آن با به‌کارگیری روش k -متوسط داده‌های کاهش‌یافته خوشه‌بندی می‌گردد. تحلیل‌های انجام شده در این تحقیق نشان می‌دهد که خوشه‌بندی داده‌های حاصل از NMF+PCA نتایج بهتری را ارائه می‌دهد.

فهرست مطالب

۱	مقدمه	۱
۹	تعاریف و پیش‌نیازها	۲
۱۱	۱.۲ مقدمه	۱.۲
۱۱	۲.۲ بیوانفورماتیک	۲.۲
۱۴	۳.۲ پیش‌نیازهای زیستی	۳.۲
۲۴	۴.۲ داده‌کاوی	۴.۲
۳۱	۵.۲ مفاهیم ریاضی	۵.۲
۳۱	۱.۵.۲ ترکیب مخروطی و محدب	۱.۵.۲
۳۵	کاهش بعد	۳
۳۷	۱.۳ مقدمه	۱.۳
۳۷	۲.۳ کاهش بعد	۲.۳
۴۰	۳.۳ تحلیل مؤلفه اصلی	۳.۳
۴۱	۱.۳.۳ الگوریتم <i>PCA</i>	۱.۳.۳
۵۰	۴.۳ تجزیه‌ی ماتریس نامنفی	۴.۳
۵۵	۱.۴.۳ کاهش بعد و استخراج ویژگی‌ها	۱.۴.۳
۵۹	۲.۴.۳ الگوریتم تجزیه‌ی ماتریس نامنفی	۲.۴.۳
۶۴	۳.۴.۳ بهینه‌سازی تجزیه‌ی ماتریس نامنفی	۳.۴.۳
۶۵	۴.۴.۳ مقداردهی اولیه	۴.۴.۳

۷۴	نتیجه‌گیری	۵.۳
۷۷	خوشه‌بندی	۴
۷۹	مقدمه	۱.۴
۷۹	یادگیری با نظارت در مقابل یادگیری بدون نظارت	۲.۴
۷۹	طبقه‌بندی	۳.۴
۸۰	خوشه‌بندی	۴.۴
۸۲	کاربردهای خوشه‌بندی	۱.۴.۴
۸۴	انواع خوشه‌ها	۲.۴.۴
۸۴	ویژگی‌های روش خوشه‌بندی خوب	۳.۴.۴
۸۵	ویژگی‌های الگوریتم‌های خوشه‌بندی	۴.۴.۴
۸۶	خوشه‌بندی در مقابل چندی‌سازی برداری	۵.۴.۴
۸۷	انواع روش‌های خوشه‌بندی	۶.۴.۴
۸۸	خوشه‌بندی k -متوسط	۵.۴
۹۰	نتیجه‌گیری	۶.۴
۹۱	جمع آوری داده‌ها و ارزیابی نتایج الگوریتم‌ها	۵
۹۳	مقدمه	۱.۵
۹۳	جمع آوری داده‌ها	۲.۵
۹۳	آماده‌سازی داده‌ها	۳.۵
۹۵	NMF در مقایسه با PCA برای کاهش داده‌های ریزآرایه	۴.۵
۹۷	خوشه‌بندی داده‌های حاصل از الگوریتم PCA	۵.۵
۹۸	کاهش بعد داده‌ها با الگوریتم NMF	۶.۵
۱۰۰	مقایسه الگوریتم‌های مختلف تجزیه‌ی ماتریس نامنفی	۷.۵
۱۰۴	نتیجه‌گیری	۸.۵
۱۰۶	مراجع	

۱۱۴

واژه‌نامه فارسی به انگلیسی

۱۱۷

واژه‌نامه انگلیسی به فارسی

لیست تصاویر

۱۷	مدل مارپیچی DNA	۱.۲
۱۸	ساخته شدن پروتئین از DNA	۲.۲
۲۰	هیبرید شدن	۳.۲
۲۲	فناوری ریزآرایه	۴.۲
۴۱	کاربرد PCA با توزیع دوبعدی	۱.۳
۴۴	نمایش داده‌های اصلی	۲.۳
۴۶	رسم بردارهای مشخصه	۳.۳
۴۷	مؤلفه‌های اصلی مجموعه داده دو بعدی	۴.۳
۵۰	داده‌های نهایی	۵.۳
۵۴	تجزیه نامنفی ماتریس به‌عنوان یک تبدیل مختصات کانونی	۶.۳
۵۷	مقایسه کیفیت تقریب سه روش کاهش بعد کلاس ماتریس T	۷.۳
۵۷	شش پازل ۲۵ پیکسلی	۸.۳
۵۸	پازل‌های پایه نهایی ($r = ۳$)	۹.۳
۵۹	خطای تقریب پازل‌های	۱۰.۳
۶۱	پیشرفت پازل پایه	۱۱.۳
۶۲	رفتار تابع هدف Θ_{NMFE}	۱۲.۳
۸۲	خوشه‌بندی حیوانات	۱.۴
۹۷	خوشه‌بندی داده‌های حاصل از روش PCA با استفاده از روش k -متوسط	۱.۵

۹۹	<i>NNDSVDar</i> و <i>NNDSVDa</i> ، <i>NNDSVD</i> اولیه بامقداردهی اولیه <i>NMF</i> مقایسه	۲.۵
۱۰۱	مقایسه <i>NMF</i> بامقداردهی اولیه مختلف	۳.۵
	خوشه‌بندی داده‌های بیان ژن حاصل از کاهش الگوریتم تجزیه‌ی ماتریس نامنفی بر روی	۴.۵
۱۰۲	پایگاه داده <i>gems – system</i>	

فصل ۱

مقدمه

مقدمه

علم زیست شناسی از قدیمی‌ترین علمی است که بشر به آن توجه داشته است؛ اما از حدود یک قرن پیش، این علم وارد مرحله جدیدی به نام ژنتیک شد. با پیشرفت علم ژنتیک، ژن به‌عنوان عامل کنترل‌کننده ویژگی و عملکرد سلول‌ها شناخته شد. در علم ژنتیک، هنگامی که یک ژن در سلول مورد استفاده قرار می‌گیرد، می‌گویند آن ژن بیان شده یا روشن است. تغییر بیان ژن^۱ سلول در اثر جهش‌های ژنتیکی و شرایط خاص موجب فعالیت غیر عادی سلول و بروز بیماری می‌شود. امروزه حجم عظیمی از مطالعات پزشکی در جهت شناسایی و درمان بیماری‌هایی است که از طریق ژن منتقل می‌شود. همچنین برای بررسی و نگهداری اطلاعات ژنتیکی، فناوری‌های مفیدی به وجود آمده‌است که یکی از این تکنولوژی‌ها، ریزآرایه^۲ می‌باشد.

در یک آزمایش ریزآرایه، توالی ژن‌ها به وسیله رنگ‌های فلورسنت متفاوت علامت‌گذاری شده و روی یک آرایه شیشه‌ای یا پلاستیکی هیبرید می‌شوند. پس از آن، آرایه مزبور توسط اشعه لیزر پوشش داده و تصاویری شامل هزاران لکه رنگی ایجاد می‌شود که انجام محاسبات عددی روی آن‌ها اندازه‌ی بیان ژن را نتیجه می‌دهد. سپس داده‌های حاصل از آزمایش در ماتریس بیان ژن ذخیره می‌شود. سطرهای این ماتریس را ژن‌ها و ستون‌های آن را افراد نمونه تشکیل می‌دهد. ماتریس مزبور در کشف و شناسایی پدیده‌های پزشکی ناشناخته، رفع اثرات ناهنجار ژنتیکی و پیشگیری و درمان بیماری‌ها مورد استفاده قرار می‌گیرد.

فناوری ریزآرایه امکان بررسی بیان هزاران ژن را به‌صورت همزمان فراهم می‌سازد و حجم عظیمی از داده‌های بیان ژنی را تولید می‌کند که می‌توان با استفاده از روش‌های داده‌کاوی^۳ دانش مورد نیاز جهت شناسایی و درمان بیماری‌ها را از آن استخراج نمود. در داده‌های ریزآرایه، تعداد ژن‌ها بسیار زیاد و تعداد نمونه‌ها کم می‌باشد. در تجزیه و تحلیل داده‌ها، وجود همه‌ی ژن‌ها ضروری نیست، یعنی ممکن است بسیاری از آن‌ها نیاز نباشند و تعدادی دیگر حاوی اطلاعات تقریباً یکسانی باشند و باعث افزونگی^۴ شوند، پس عمل ساده‌سازی و انتخاب

^۱Gene expression

^۲Microarray

^۳Data mining

^۴Redundancy

ژن‌های مناسب امری ضروری است که این عمل به کمک روش‌های کاهش بعد^۵ انجام می‌شود.

داده‌کاوی یکی از مراحل کشف دانش^۶ است. در کشف دانش ابتدا داده‌های مناسب با توجه به نوع کاربرد جمع‌آوری و سپس پیش‌پردازش می‌شود بدان معنا که مقادیر گم‌شده^۷ و پرت^۸ آن‌ها جایگزین شده، تبدیل و نرمال‌سازی روی آن‌ها انجام می‌شود. در جایی که حجم مجموعه داده بزرگ باشد عمل کاهش بعد داده‌ها نیز انجام می‌شود. هدف در این عمل، تولید مجموعه‌ای کوچک‌تر از داده‌های اولیه است به گونه‌ای که نتایج حاصل از این مجموعه با نتایج حاصل از داده اولیه تقریباً یکسان باشد. پس از اجرای مراحل پیش‌پردازش و کاهش بعد داده، فرآیند داده‌کاوی انجام می‌شود. با استفاده از روش‌های داده‌کاوی می‌توان دانش مفید را از داده‌ها استخراج نمود. از روش‌های کاهش بعد می‌توان برای شناسایی ابعاد و ویژگی‌های مطلوب برای نمایش داده‌ها در فرآیند داده‌کاوی استفاده نمود.

تاکنون روش‌های زیادی از جمله روش‌های تحلیل مؤلفه اصلی^۹، تحلیل مؤلفه مستقل^{۱۰}، تجزیه‌ی نیمه‌گسسته^{۱۱}، تجزیه‌ی مقادیر تکین^{۱۲} و تجزیه‌ی ماتریس نامنفی^{۱۳} برای کاهش بعد داده‌ها معرفی شده است. در برخی از کاربردها به‌عنوان مثال تحلیل و داده‌کاوی داده‌های ریزآرایه، لازم است به منظور کاهش بعد از روش‌هایی استفاده نمود که شرایط نامنفی بودن ماتریس‌های کاهش یافته را حفظ کند که از بین تمام روش‌های مزبور تنها روش تجزیه‌ی ماتریس نامنفی این شرایط را حفظ می‌کند.

تجزیه‌ی ماتریس نامنفی برای اولین بار در سال ۱۹۷۰ تحت عنوان «تجزیه‌ی منحنی‌های خود سازمانده»^{۱۴} ارائه شد. این تجزیه با نام تجزیه‌ی ماتریس مثبت به صورت گسترده توسط محققان اسپانیایی مطالعه گردید.

^۵Dimension reduction

^۶Knowledge discovery

^۷Missing values

^۸Outlier

^۹Principal component analysis(PCA)

^{۱۰}Independent component analysis(ICA)

^{۱۱}Semi-discrete decomposition(SSD)

^{۱۲}Singular value decomposition(SVD)

^{۱۳}Nonnegative matrix factorization(NMF)

^{۱۴}Self modeling curve resolution

در سال ۱۹۹۹ تحقیقات «لی» و «سونگ»^{۱۵} موجب توجه محققان به این تجزیه و کاربرد آن در زمینه داده کاوی و یادگیری ماشین شد. پس از انتشار مقاله آن‌ها [۴۲]، تجزیه‌ی مزبور با نام تجزیه‌ی ماتریس نامنفی معرفی گردید. این تجزیه بسیاری از ساختارهای داده‌ی اصلی را حفظ و نامنفی بودن هر دو ماتریس پایه و ضرایب را ضمانت می‌کند. در تحلیل داده‌های بیان ژن ضروری است که پایه‌ها نامنفی نگه داشته شود، بنابراین برای کاهش بعد این داده‌ها از روش تجزیه‌ی ماتریس نامنفی استفاده شده است.

پس از کاهش بعد داده‌ها از روش‌های داده کاوی جهت استخراج الگوهای مفید از داده‌ها استفاده می‌شود. خوشه‌بندی^{۱۶} یکی از روش‌های داده کاوی است که به مفهوم تقسیم داده‌ها به گروه‌هایی از داده‌های مشابه می‌باشد. قابل ذکر است که اعتبار نتایج این روش به انتخاب معیار شباهت بین داده‌ها وابسته است. انتخاب معیار شباهت مناسب، در داده‌هایی که دارای تعدد ویژگی یا بعد زیاد می‌باشد، کاری مشکل است و ممکن است به نتایج نامطلوب منجر شود. همچنین پیچیدگی محاسباتی روش‌های خوشه‌بندی زیاد بوده و در برخورد با مجموعه داده‌های بزرگ ناپایدار است. جهت رفع این مشکل‌ها می‌توان از روش‌های کاهش بعد در جهت تولید مجموعه داده‌هایی با حجم کوچک‌تر و ویژگی‌های مطلوب از مجموعه داده‌های اصلی استفاده نمود و آن را به‌عنوان ورودی روش خوشه‌بندی به کار برد. روش‌های خوشه‌بندی به دو دسته کلی خوشه‌بندی سلسله مراتبی^{۱۷} و خوشه‌بندی افزایی^{۱۸} تقسیم می‌شوند. روش خوشه‌بندی k -متوسط^{۱۹} یکی از روش‌های خوشه‌بندی افزایی است. الگوریتم این روش بسیار ساده و قابل فهم است و عملکرد آن شبیه عملکرد روش تجزیه‌ی ماتریس نامنفی است و همچنین دارای نتایج مطلوب در زمینه‌ی خوشه‌بندی داده‌ها می‌باشد [۳۲]، [۵۱].

در این تحقیق روش تجزیه‌ی ماتریس نامنفی به‌عنوان یک روش تکراری در کاهش بعد داده‌های بیان ژن استفاده خواهد شد. ایده اصلی در یک روش تکراری این است که ابتدا مسئله به یک شکل معادل نوشته شود و سپس با شروع از یک جواب اولیه، دنباله‌ای همگرا به جواب اصلی تولید گردد. در روش تجزیه‌ی ماتریس نامنفی می‌توان از روش‌های مختلفی برای مقداردهی اولیه استفاده نمود. الگوریتمی که توسط «لی» و «سونگ» ارائه شد از روش تصادفی جهت مقداردهی اولیه عوامل تجزیه استفاده می‌نمود. تاکنون مطالعات زیادی جهت

^{۱۵}Lee and Seung

^{۱۶}Clustering

^{۱۷}Hierarchical clustering

^{۱۸}Partition clustering

^{۱۹}k-means

استفاده از یک روش مقداردهی اولیه مناسب برای تجزیه NMF انجام گرفته است. نتایج این مطالعات نشان داده است که روش تجزیه مقدار تکین مضاعف نامنفی^{۲۰} نسبت به روش‌های دیگر دارای نتایج مطلوب‌تری است. در این رساله، تجزیه ماتریس نامنفی با مقداردهی اولیه تصادفی و تجزیه مقدار تکین مضاعف نامنفی مطالعه و روش تحلیل مؤلفه اصلی به منظور مقداردهی اولیه تجزیه NMF پیشنهاد و بررسی خواهد شد. پس از کاهش بعد داده‌های ریزآرایه، از روش k -متوسط برای خوشه‌بندی داده‌ها استفاده می‌شود.

اهمیت و ضرورت تحقیق

امروزه بیماری‌هایی که از طریق ژن منتقل می‌شوند و از جمله آن‌ها انواع سرطان‌ها و بیماری‌های قلبی، یکی از عوامل مهم مرگ و میر است. به‌همین دلیل حجم عظیمی از مطالعات پزشکی جهت شناسایی، درمان و پیشگیری از این بیماری‌ها می‌باشد و برای این منظور ضروری است که عوامل بیماری شناسایی و طبقه‌بندی شود تا مراحل درمانی مناسب با نوع خاص بیماری انجام شود. فناوری ریزآرایه و روش‌های داده‌کاوی کمک می‌کند تا با تحلیل تغییرات بیان هزاران ژن، به صورت همزمان بتوان فرایند ایجاد بیماری‌هایی مثل سرطان را شناسایی کرد و در درمان این بیماری‌ها گام‌های مهمی برداشت.

یکی از عوامل بسیار مهم سرطان، جهش‌های ژنتیکی است. این جهش‌ها بسیار شبیه به هم هستند و نیاز به آزمایش‌های دقیق جهت شناسایی آن‌ها از یکدیگر وجود می‌باشد. با شناسایی کامل ژنوم انسانی و ژن‌هایی که در این جهش‌ها موثر هستند، انتظار می‌رود تشخیص بین این جهش‌ها با استفاده از روش‌هایی بر مبنای بیان ژن، امکان‌پذیر باشد. محققان تلاش می‌کنند تا با خوشه‌بندی بیماری‌های ژنتیکی بر اساس رفتارها و عوامل ژنتیکی در زیر گروه‌های همگن، فرایند تشخیص و درمان بیماری‌ها را سرعت بخشند.

اهداف اصلی تحقیق

اهداف اصلی این پژوهش عبارتند از:

۱. مطالعه‌ی تجزیه‌ی ماتریس نامنفی

۲. مطالعه و بررسی روش تحلیل مؤلفه اصلی جهت مقداردهی اولیه‌ی تجزیه‌ی ماتریس نامنفی

^{۲۰}Nonnegative double singular value decomposition(NNDSVD)

۳. مطالعه و بررسی روش تجزیه‌ی مقدار تکین مضاعف نامنفی جهت مقداردهی اولیه‌ی تجزیه‌ی ماتریس

نامنفی

۴. تحلیل نتایج خوشه‌بندی داده‌های حاصل از روش تجزیه‌ی ماتریس نامنفی

ساختار تحقیق

در این تحقیق، فصل اول مقدمات و کلیات موضوع را ارائه داده است. فصل دوم شامل تعاریف و پیش زمینه‌های لازم از جمله بیوانفورماتیک، پیش زمینه‌های زیستی، مفاهیم داده‌کاوی و تعاریف ریاضی می‌باشد. فصل سوم روش‌های تحلیل مؤلفه اصلی و تجزیه‌ی ماتریس نامنفی را جهت کاهش بعد داده‌ها و روش تجزیه‌ی مقدار تکین مضاعف نامنفی را جهت مقداردهی اولیه‌ی روش تجزیه‌ی ماتریس نامنفی شرح می‌دهد. در فصل چهارم خوشه‌بندی و روش‌های آن توضیح و روش خوشه‌بندی k -متوسط از بین این روش‌های شرح داده می‌شود. فصل پنجم ابتدا به آماده‌سازی و پیش‌پردازش داده‌ها پرداخته است و سپس نتایج کاهش بعد داده‌ها را با استفاده از هر یک از الگوریتم‌های معرفی شده در پایان‌نامه مورد بررسی و مقایسه قرار می‌دهد و نتایج خوشه‌بندی داده‌های بیان ژن یکی از پایگاه‌داده‌ها را مورد تحلیل قرار می‌دهد.

فصل ۲

تعاریف و پیش‌نیازها

