

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشکده‌ی مهندسی برق و کامپیوتر

پایان‌نامه‌ی کارشناسی ارشد در رشته‌ی مهندسی کامپیوتر - هوش مصنوعی

ارائه‌ی یک مدل آماری برای خوشه‌بندی متون

به وسیله‌ی
حمید محمودی

استاد راهنما
دکتر اقبال منصوری

تیر ماه ۱۳۹۱

به نام خدا

اظهار نامه

اینجانب حمید محمودی دانشجوی رشته ی مهندسی کامپیوتر گرایش هوش مصنوعی، دانشکده مهندسی برق و کامپیوتر اظهار می کنم که این پایان نامه حاصل پژوهش خودم بوده و در جاهایی از منابع دیگران استفاده کرده ام.نشانی دقیق و مشخصات کامل آنرا نوشته ام.همچنین اظهار می کنم که تحقیق و موضوع پایان نامه ام تکراری نیست و تعهد می نمایم که بدون مجوز دانشگاه دستاورد های آنرا منتشر ننموده و در اختیار غیر قرار ندهم.کلیه حقوق این اثر مطابق با آیین نامه مالکیت فردی و معنوی متعلق به دانشگاه شیراز است.

نام و نام خانوادگی: حمید محمودی

تاریخ و امضا: ۱۳۹۱/۵/۱۴

بنام خدا

ارائه‌ی یک مدل آماری برای خوشه‌بندی متون

به کوشش

حمید محمودی

پایان نامه

ارائه شده به تحصیلات تکمیلی دانشگاه به عنوان بخشی
از فعالیت‌های تحصیلی لازم برای اخذ درجه کارشناسی ارشد

در رشته

مهندسی کامپیوتر-هوش مصنوعی

از دانشگاه شیراز

شیراز

جمهوری اسلامی ایران

ارزیابی شده توسط کمیته پایان نامه یا درجه: عالی

..... دکتر اقبال منصوری، استادیار بخش مهندسی و علوم کامپیوتر (رئیس کمیته)

..... دکتر فرشاد تاجری پور، استادیار بخش مهندسی و علوم کامپیوتر

..... دکتر محمدهادی صدرالدینی، دانشیار بخش مهندسی و علوم کامپیوتر

مرداد ۱۳۹۱

سپاسگزاری

برخود لازم می‌دانم از زحمات استاد ارجمند جناب آقای دکتر اقبال منصوری که همواره در تمام مدت انجام پایان‌نامه مرا از راهنمایی‌ها و مساعدت‌های بی‌دریغ‌شان بهره‌مند نمودند تشکر کنم. همچنین از استادان مشاور خود جناب آقای دکتر فرشاد تاجری‌پور و جناب آقای دکتر محمدهادی صدرالدینی که با نظرات و راهنمایی‌های مفیدشان مرا در پیشبرد این پایان‌نامه یاری نموده‌اند کمال تشکر را دارم.

در انتها از تمام عزیزانی که مرا در انجام این پروژه تحقیقاتی یاری نمودند کمال تشکر و قدردانی را دارم.

چکیده

ارائه‌ی مدل آماری برای خوشه‌بندی متون

به کوشش

حمید محمودی

در این پایان‌نامه سعی در ارائه‌ی یک مدل آماری برای خوشه‌بندی متون داشته‌ایم. هر خوشه به منزله‌ی جزئی از یک مدل ترکیبی در نظر گرفته می‌شود که شامل پارامترهای عدد اولویت جزء، بردار میانگین و ماتریس کواریانس جزء می‌باشد. هدف از ارائه‌ی یک مدل آماری، در نظر گرفتن پخشش‌های متفاوت برای مجموعه داده‌هایی که لزوماً پخشش آنها کروی نیست، می‌باشد. الگوریتمهای خوشه‌بندی مانند *K-Means* و مشتقات آن که با یک پارامتر-که معمولاً مراکز خوشه‌هاست- کار می‌کنند، سعی دارند خوشه‌هایی با پخشش کروی را ایجاد کنند که این در مورد همه‌ی مجموعه داده‌های دنیای واقعی صدق نمی‌کند. هدف دیگر این پایان‌نامه، به مقدار دهی اولیه‌ی پارامترهای مدل بر می‌گردد. بسیاری از کارهای انجام شده در زمینه‌ی خوشه‌بندی بدون نظارت متون، ساخت چندین مدل مختلف با مقداردهی‌های اولیه‌ی تصادفی بوده است و نهایتاً مدل برتر را بر اساس یک معیار خاص انتخاب می‌کردند. از آنجائیکه مقداردهی تصادفی در همه‌ی اجراها نتایج قابل اعتماد و منظمی ارائه نمی‌دهد، ما با ارائه‌ی یک رویه‌ی چند مرحله‌ای و بدون کمک از هرگونه ناظر خارجی و با استفاده از الگوریتمهای سلسله مراتبی که پیچیدگی محاسباتی آن‌را با انتخاب یک مجموعه کوچک از نمونه‌های برتر و همچنین کاهش فضای ابعاد، کاهش دادیم موفق شدیم برای هر خوشه بهترین متون مربوطه را بیابیم و بوسیله‌ی آنها پارامترهای مذکور را مقداردهی کنیم. نتایج آزمایشگاهی و نمودارهای مقایسه‌ای به صورت کاملاً واضح نشان می‌دهند که روش‌های پیشنهاد شده در این پایان‌نامه از عملکرد بالاتری نسبت به روشهای ارائه شده‌ی مشابه داشته‌اند.

واژگان کلیدی: خوشه‌بندی بدون نظارت متون، روش بیشینه کردن امید، خوشه‌بندی سلسله مراتبی.

فهرست مطالب

صفحه	عنوان
	فصل اول: مقدمه و تعاریف اولیه
۲	۱-۱- سیستم‌های بازیابی اطلاعات
۳	۱-۲- خوشه‌بندی متون
۳	۱-۳- اهمیت خوشه‌بندی متون
۵	۱-۴- انگیزه انجام این پایان‌نامه
۶	۱-۵- پیش‌پردازش متون
۶	۱-۵-۱- نشانه‌گذاری
۶	۱-۵-۲- حذف کلمات بی‌اثر
۶	۱-۵-۳- ریشه‌یابی
۷	۱-۶- مدل فضا برداری
۷	۱-۶-۱- مدل دودویی
۸	۱-۶-۲- مدل وزن‌دهی عبارت
۱۰	۱-۶-۳- محاسن و معایب مدل فضا برداری
۱۱	۱-۷- معیارهای شباهت
۱۱	۱-۷-۱- فاصله اقلیدسی
۱۱	۱-۷-۲- معیار شباهت کسینوسی
۱۳	۱-۷-۳- معیار شباهت Jaccard
۱۳	۱-۷-۴- تابع واگرایی Jensen Shanon (JS)
۱۴	۱-۷-۵- تابع واگرایی Kulback-Leibler
۱۴	۱-۸- روشهای خوشه‌بندی
۱۵	۱-۸-۱- خوشه‌بندی بر اساس تقسیم‌بندی متون
۲۶	۲-۸-۲- الگوریتمهای خوشه‌بندی سلسله‌مراتبی

فصل دوم: انتخاب ویژگیها با استفاده از معیارهای آماری

۳۵	۲-۱- مقدمه
۳۶	۲-۱-۱- روشهای فیلتر
۳۶	۲-۱-۲- روشهای پوشاننده
۳۷	۲-۲- روشهای معمول انتخاب ویژگی فیلتر برای یادگیری بدون ناظر
۳۸	۲-۲-۱- فرکانس متون (DF)
۳۸	۲-۲-۲- روش استحکام کلمه (TS)
۳۸	۲-۲-۳- رتبه بندی بر اساس آنتروپی
۳۹	۲-۲-۴- مشارکت کلمه (TC)
۴۰	۲-۲-۵- انتخاب ویژگی بر اساس χ^2
۴۱	۲-۲-۶- بهره اطلاعات (IG)

فصل سوم: روشهای مختلف ارزیابی خوشه بندی

۴۴	۳-۱- ارزیابی صحت اعتبار خوشه ها
۴۶	۳-۲- روشهای مختلف اعتبارسنجی
۴۷	۳-۲-۱- اعتبارسنجی داخلی
۵۱	۳-۲-۲- اعتبارسنجی خارجی

فصل چهارم: ارائه راه حل پیشنهادی برای مسئله ی خوشه بندی متون

۵۷	۴-۱- انگیزه از ارائه روش پیشنهادی
۶۰	۴-۲- ارائه روش تلفیقی انتخاب ویژگی
۶۲	۴-۳- خوشه بندی متون بر اساس خوشه بندی نیمه نظارتی کلمات
۶۴	۴-۴- ارائه ی یک مدل آماری بدون نظارت برای خوشه بندی متون

فصل پنجم: نتایج آزمایشات

۷۰	۵-۱- مجموعه های داده
۷۰	۵-۱-۱- Reuters-21587
۷۲	۵-۱-۲- 20Newsgroup
۷۳	۵-۲- نتایج آزمایشات

۷۳ ۵-۲-۱- نتایج الگوریتم پیشنهادی اول

۷۹ ۵-۲-۲- نتایج الگوریتم پیشنهادی دوم

فصل ششم: نتیجه‌گیری و پیشنهادات

۸۹ ۶-۱- نتیجه‌گیری

۹۰ ۶-۲- پیشنهاد برای انجام کارهای آتی

۹۱ فهرست منابع و مأخذ

فهرست جداول

صفحه	عنوان
۳۳	جدول ۱: شبه کد الگوریتم Bisecting K-Means
۷۳	جدول ۲: کلاسه‌های مجموعه داده‌ی 20newsgroup
۷۴	جدول ۳: خلاصه اطلاعات سه مجموعه داده
۷۵	جدول ۴: نتایج مقدار F سراسری روی مجموعه داده‌ی NG20
۷۵	جدول ۵: نتایج مقدار F سراسری روی مجموعه داده‌ی NG10
۷۶	جدول ۶: نتایج F سراسری روی re0 و re1
۷۹	جدول ۷: ده زیر مجموعه از مجموعه داده اصلی 20-newsgroup
۷۹	جدول ۸: سه مجموعه داده‌ی اصلی
۸۲	جدول ۹: نتایج ریزمیانگین F مربوط به ده زیر مجموعه کوچک روی هشت الگوریتم
۸۳	جدول ۱۰: نتایج ریزمیانگین مربوط به سه مجموعه داده‌ی مرجع

۸۳

فهرست اشکال

عنوان	صفحه
شکل ۱: مدل برداری بولیین از سه متن و یک پرس و جوی کاربر	۷
شکل ۲: مدل وزن‌دهی عبار	۸
شکل ۳: فاصله زاویه‌ای بین بردار q و دو بردار متن d_1 و d_2	۱۲
شکل ۴: الگوریتم استاندارد batch K-Means	۱۶
شکل ۵: نحوه خوشه‌بندی K-Means	۱۷
شکل ۶: نحوه خوشه‌بندی الگوریتم پیشینه امید در ۶ مرحله مختلف	۲۲
شکل ۷: فلوجارت الگوریتم پیشینه امید	۲۵
شکل ۸: ساختار درختی دودویی	۲۷
شکل ۹: شبه کد روشهای تجمعی	۲۸
شکل ۱۰: پدیده زنجیره ای در پیوند یگانه	۲۹
شکل ۱۱: نقاط دورافتاده دو خوشه نامشابه را در پیوند کامل ترکیب می‌کنند	۲۹
شکل ۱۲: روش میانگین گروهی	۳۰
شکل ۱۳: رویه روش مرکزی برای خوشه‌بندی ۶ متن	۳۱
شکل ۱۴: شبه کد الگوریتمهای تقسیمی سلسله مراتبی	۳۲
شکل ۱۵: مراحل اساسی خوشه‌بندی	۴۳
شکل ۱۶: (a) داده اصلی با ۳ خوشه (b) نتیجه K-means با ۴ خوشه	۴۵
شکل ۱۷: چندین مجموعه داده با پخشش‌های متفاوت	۵۹
شکل ۱۸: فلوجارت کلی الگوریتم ارائه شده	۶۴
شکل ۱۹: قطع ادغام خوشه‌ها در مرحله‌ای خاص	۶۷
شکل ۲۰: شبه کد خوشه بندی سلسله مراتبی متون	۶۷
شکل ۲۱: تقسیم‌بندی Mod-Apte Split مجموعه متون Reuters-21587	۷۱
شکل ۲۲: یک نمونه متن از کلاس earn	۷۲

شکل ۲۳: مقایسه‌ی مقادیر NMI از الگوریتمهای مختلف با کسرهای مختلف نمونه های برچسب دار روی مجموعه داده‌ی کامل 20-Newsgroup	۷۷
شکل ۲۴: مقایسه‌ی مقادیر NMI از الگوریتمهای مختلف با کسرهای مختلف نمونه های برچسب دار روی مجموعه داده‌ی 10-Newsgroup	۷۷
شکل ۲۵: مقایسه‌ی مقادیر NMI از الگوریتمهای مختلف با کسرهای مختلف نمونه های برچسب دار روی re0	۷۸
شکل ۲۶: مقایسه‌ی مقادیر NMI از الگوریتمهای مختلف با کسرهای مختلف نمونه های برچسب دار روی re1	۷۸
شکل ۲۷: منحنی ایده‌آل بین صحت-یادآوری	۸۴
شکل ۲۸: منحنی صحت-یادآوری مربوط به مجموعه داده‌ی NG10	۸۵
شکل ۲۹: منحنی صحت-یادآوری مربوط به مجموعه داده NG20	۸۶
شکل ۳۰: منحنی صحت-یادآوری مربوط به مجموعه داده‌ی Reuters(10)	۸۷

فصل اول

مقدمه

۱-۱- سیستم‌های بازیابی اطلاعات

با پیشرفت رایانه‌ها هم از جهت سخت افزار و هم از جهت نرم افزار خصوصا در دو دهه اخیر، تکنیکهای مدیریت داده‌ها و اطلاعات بیش از پیش ماشینی و خودکار شد. سیستم‌های مدیریت پایگاه داده مخصوص ذخیره و پردازش داده‌های ساختاریافته بود که مکانیزمهای جستجو برای بازیابی و پردازش اینگونه داده‌ها سراسر است و آسان بود. به علت ماهیت ویژه و غیرساختاریافته داده‌های متنی و پیچیدگی مکانیزمهای جستجو و پردازش آنها، نوع خاصی از تکنولوژی پردازش داده تحت عنوان سیستمهای بازیابی متون بوجود آمد که تفسیر کلی آن عبارت بود از تکنیکهای ذخیره و دستکاری حجمی از داده‌های متنی برای اهداف مختلف بازیابی اطلاعات.

اکثر سیستمهای بازیابی اطلاعات سنتی اطلاعات مربوط به فهرست داده‌ها را ذخیره و پردازش می‌کردند و اساس مکانیزم جستجوی آنها بر روی کلیدهای فهرست برای بازیابی منبع و آدرس محل ذخیره اطلاعات اصلی استوار بود. امروزه سیستم‌های خودکار بازیابی اطلاعات قادر هستند که هر دو موجودیت فهرست داده‌ها و کل داده‌های اصلی را پردازش و اداره کنند. بازیابی اطلاعات به مطالعه تکنیک‌های جستجو و بازیابی اطلاعات برای داده‌های غیرساختیافته متون، اطلاعاتی که درون متون هستند و ابرداده‌هایی در مورد متون اطلاق می‌شود [۱،۲]. واژه داده غیرساختیافته معمولاً به داده‌هایی اشاره می‌کند که ساختار منظم و منسجم و واضحی نداشته باشند. در مقابل داده‌های ساختیافته مانند پایگاه داده رابطه‌ای که شامل جداولی از رکوردهای مرتب هستند، دارای ساختاری منظم و ارتباطات معنایی و قابل فهم هستند. هرچند داده‌های غیرساختیافته‌ای مانند تصاویر گرافیکی، فایل‌های صوتی، فایل‌های صوتی تصویری و... نیز وجود دارند، اما کاوش‌های بازیابی اطلاعات روی بازیابی متون زبان طبیعی متمرکز است. به علت اینکه بسیاری از متون شامل سرصفحه‌های ساختار یافته - که شامل ابرداده‌هایی از قبیل عنوان، توضیحات کوتاه در مورد متن، موضوع متن، نویسنده، ناشر، تاریخ‌ها و... - و بدنه غیرساختیافته هستند، گاهی از متون به عنوان داده‌های نیمه ساختیافته نیز یاد می‌شود.

۲-۱- خوشه بندی متون

خوشه بندی یکی از پرکاربردترین کارها در حوزه داده کاوی است، که به منظور تعیین توزیع-های نمونه‌هایی که از آن توزیع تبعیت می‌کنند و همچنین کشف گروه‌هایی که تا حد ممکن از هم متفاوت و تا حد ممکن نمونه‌های درون گروهی مشابه باشند استفاده می‌شود [۳,۴]. خوشه بندی معمولاً به دسته بندی مجموعه ای از متون در خوشه های مختلف بر اساس میزان شباهتی که بین آنها وجود دارد اطلاق می‌شود. به طور ایده آل تا جایی که ممکن است خوشه ها از هم جدا و شباهت درون خوشه ای حداکثر می‌باشد. اما معمولاً همپوشانی خوشه-ها، وجود داده های نامرتبط با داده های دیگر و وجود ویژگیهای نویزی که اثر منفی روی جداکنندگی داده ها از هم دارند، سه چالش اصلی در خوشه بندی متون و دیگر حوزه های بازیابی اطلاعات هستند.

۳-۱- اهمیت خوشه‌بندی متون

در ابتدا خوشه‌بندی متون برای بهبود دقت و یادآوری یک سیستم بازیابی اطلاعات و یک راه موثر برای یافتن نزدیکترین همسایه متون مورد استفاده قرار می‌گرفت. امروزه از خوشه‌بندی برای تنظیم و سازماندهی نتایج موتورهای جستجوگر برای پاسخ به پرس‌وجوهای کاربران مورد استفاده واقع می‌شود. یکی دیگر از کاربردهای مهم خوشه بندی طبقه‌بندی خودکار مجموعه-های بزرگی از متون به صورت سلسله مراتبی می‌باشد.

یکی از مهمترین انگیزه‌های استفاده از خوشه‌بندی، تعیین و آشکار کردن ساختار ذاتی و پنهان یک مجموعه داده است. کاربران انسانی به علت تفاوت در سلیقه و طرز تفکرات مختلف از کشف ساختار ذاتی و درونی مجموعه داده‌ای بزرگ متون ناتوان هستند. به عنوان مثال فرض کنید مجموعه‌ای از متون مربوط به بازیهای المپیک را در اختیار داریم. یک متخصص انسانی برای خوشه‌بندی ممکن است متون را بر اساس نام بازیهای مختلف (فوتبال، هاکی، شنا،...)، متخصص دیگری بر اساس نام کشورهای شرکت‌کننده (انگلیس، استرالیا،...) و دیگری بر اساس نوع بازیها (زمستانی یا غیر زمستانی، با توپ یا بدون توپ، قدرتی یا سرعتی،...) خوشه‌بندی کنند و در نهایت چندین خوشه‌بندی متفاوت داشته باشیم و به علت اینکه هیچ یک نتوانسته‌اند بر اساس ساختار درونی خوشه‌بندی کنند نمی‌توانیم هیچ یک را بر دیگری برتری دهیم.

یکی دیگر از عوامل مهم استفاده از خوشه‌بندی، گسترش روزافزون اطلاعات خصوصاً در حوزه وب و اینترنت می‌باشد که خوشه‌بندی و طبقه‌بندی اطلاعات را جزء جداناپذیر

موتورهای جستجوگر مدرن امروزی کرده است. موتورهای جستجوگر امروزی برای یافتن بهترین نتایج مربوط به یک پرس‌وجو به جای جستجو در میان میلیونها صفحه وب که از نقطه نظر زمانی برای یک کاربر معمولی قابل قبول نیست، کوشش می‌کند تا در میان خوشه‌های مختلف که تعدادشان به مراتب کمتر از تعداد صفحات وب هستند جستجو کنند.

کاربردهای خوشه‌بندی

کاهش داده^۱: آنالیز خوشه‌بندی می‌تواند در فشرده‌سازی و یکجا کردن اطلاعات موجود در داده‌ها به کار رود. در بسیاری از مواقع مجموعه داده‌ای در اختیار، آنقدر وسیع و گسترده است که نیاز به طبقه‌بندی و فشرده‌گی تا حد ممکن داده‌ها دارد.

تولید فرضیه^۲: این مورد برای استنتاج یا کشف فرضیه‌هایی که پنهان در داده‌ها هستند کاربرد دارد. برای مثال پایگاه داده‌ی خرده‌فروشی را در نظر بگیرید که دو گروه از مشتریها را بر اساس سن افراد و زمان خرید دارد. می‌توانیم به عنوان مثال یک گمانه کلی در این مورد داشته باشیم که افراد جوان در عصر یا افراد مسن در صبح به خرید می‌روند.

پیش‌بینی بر اساس خصوصیات خوشه‌ها: خوشه‌بندی مجموعه‌ای از داده‌ها را به گروه‌های متفاوتی بر اساس ویژگی‌هایشان تقسیم می‌کند. داده‌هایی که پس از خوشه‌بندی به دست می‌آید و در فرایند خوشه‌بندی دخالت نداشته‌اند، می‌توانند بر اساس شباهتی که با نزدیکترین خوشه را دارند به آن خوشه تخصیص یابند. به عنوان مثال فرض کنید داده‌هایی در مورد بیمارانی داریم که با یک الگوریتم خوشه‌بندی بر اساس واکنشی که به داروهای مختلفی از خود بروز می‌دهند، خوشه‌بندی می‌شوند. پس برای بیمار جدیدی، و بر اساس نزدیکترین خوشه‌ای که به آن تخصیص می‌یابد می‌توان معالجه‌ای که روی بیماران آن خوشه انجام می‌شود، برای او نیز تجویز کرد.

تجارت: در تجارت، خوشه‌بندی می‌تواند به بازاریها در کشف گروه‌های معنی داری از مشتری‌هایشان بر اساس خریدهایی که می‌کنند کمک کند.

زیست‌شناسی: در زیست‌شناسی می‌توان از خوشه‌بندی در طبقه‌بندی و دسته‌بندی ژنها بر اساس عاملیت^۳ آنها در ساختارهای موروثی در جمعیت‌ها استفاده کرد.

دیگر کاربردهای خوشه‌بندی می‌توان به وب‌کاوی، متن‌کاوی، پردازش تصویر، تصاویر ماهواره‌ای، تجهیزات پزشکی، سیستم اطلاعاتی جغرافیایی^۴ (GIS) اشاره کرد.

¹ Data reduction

² Hypothesis Generation

³ Functionality

⁴ Geographical Information System

۴-۱- انگیزه‌ی انجام این پایان‌نامه

بسیاری از الگوریتمهای ارائه شده در حوزه‌ی خوشه‌بندی متون اهمیتی به پخشش داده‌ها نمی‌دهند. الگوریتمهای مشهوری مانند *K-Means* و مشتقات آن مانند *K-Means* کروی، *K-Means* فازی،... پخشش داده‌ها را به صورت کروی یا ابرکروی در نظر می‌گیرند. علت این موضوع به متریک‌های معیار فاصله یا عدم شباهت اقلیدسی یا کسینوسی برمی‌گردد. ما در این پایان‌نامه قصد داریم تا با ارائه یک مدل احتمالی گوسی از توزیع خوشه‌ها در فضای بزرگ متن-عبارت، ضمن در نظر گرفتن پخشش داده‌ها در ابعاد مختلف بتوانیم مجموعه‌هایی با توزیع غیرکروی را مدل کنیم.

یکی دیگر از اهداف ما، توجه به مقداردهی اولیه برای پارامترهای خوشه‌ها قبل از عمل خوشه‌بندی است. مقداردهی اولیه در *K-Means* شامل تعیین تنها پارامتر، یعنی مراکز اولیه خوشه‌های موجود در مجموعه است. انتخاب تصادفی پارامترها منجر به غیرقابل اعتماد و نوسانی بودن نتایج در اجراهای مختلف الگوریتم می‌شود.

ما در این پایان‌نامه یک روش جدید انتخاب ویژگی را ارائه نموده ایم که ضمن اینکه به میزان زیادی ابعاد مسئله را کاهش می‌دهد، می‌تواند بهترین و ارزشمندترین ویژگی‌ها (ابعاد) را نیز انتخاب کنید.

ساختار این پایان‌نامه بر این اساس است؛ در ادامه این فصل، به بررسی مفاهیم اولیه پرداخته شده می‌پردازیم. در فصل دوم به پیاده‌سازی و بررسی روش‌های موجود در این حوزه پرداخته می‌شود و نقاط قوت و ضعف هر روش را بیان می‌کنیم. در فصل سوم به روش‌های مختلف ارزیابی و اعتبارسنجی نتایج خوشه‌بندی ارائه شده توسط محققین مختلف خواهیم پرداخت. ارائه روش پیشنهادی و امکان بهبود کارایی آن نسبت به کارهای انجام شده همراه با جزئیات پیاده‌سازی، در فصل چهارم مطرح می‌شود. در فصل پنجم نتایج تجربی و نتایج آزمایشگاهی را که برآمد کار این پایان‌نامه است، با تفصیل مورد ارزیابی و تحلیل قرار داده می‌شود. در پایان، نتیجه‌گیری کلی از کار انجام شده را توضیح می‌دهیم و چندین پیشنهاد را برای کارهای آینده معرفی می‌کنیم.

۵-۱- پیش پردازش متون

قبل از اعمال پردازشهایی از قبیل انتخاب ویژگی و عمل خوشه بندی نیاز است که پردازشهای مقدماتی روی متون صورت گیرد تا تفسیری قابل فهم برای ماشین تولید شود. خروجی این مرحله اهمیت زیادی برای افزایش صحت نتایج و کاهش مرتبه زمانی دارد.

۱-۵-۱- نشانه گذاری^۱

در اولین مرحله پیش پردازش، می بایست متن را به صورت عباراتی جدا از هم تقسیم کرد. هر عبارت ممکن است فقط یک کلمه یا چندین کلمه که یک اصطلاح را تشکیل می دهند باشد، ولی ما همه عبارات را تک کلمه ای در نظر می گیریم و فرض می کنیم عبارت ها از هم مستقل هستند و ترتیب، اثری در آنها ندارد.

۲-۵-۱- حذف کلمات بی اثر^۲

تعداد این کلمات در زبان انگلیسی معمولاً کمتر از ۶۰۰ عدد است و تعداد آنها برای لیست های مختلف متغیر است. در این رساله یکی از لیست های استاندارد که شامل ۵۷۱ کلمه بی اثر^۳ (اضافه) است استفاده شده است. معمولاً هر متنی تعدادی کلمات موجود در این لیست را داراست که باید حذف شوند. این کلمات اثر جداکنندگی ندارند و اغلب به عنوان نویز در نظر گرفته می شوند. مثالی از این کلمات عبارتند از ("a", "the", "if", ...).

۳-۵-۱- ریشه یابی^۴

در زبان انگلیسی، معمولاً کلمات هم خانواده در چندین حرف در انتهای کلمات باهم اختلاف دارند. همه آنها یک ریشه دارند و به علت اینکه مطابقت دو کلمه به صورت مطابقت نظیر به نظیر تک تک حروف انجام می شوند، لازم است تا کلمات هم خانواده به یک شکل ریشه شان کاهش یابند. رویه ی استاندارد که برای حذف پسوند و پیشوند در این رساله استفاده شده است، Porter-Stemmer نام دارد [۶]. این مرحله پیش پردازش باعث می شود کلمات به تعداد قابل توجهی کاهش یابد و صحت نتایج را نیز افزایش دهد.

¹ Tokenization

² Stop-Word

³ <ftp://cs.cornell.edu/pub/smart/english.stop>

⁴ Stemming

۶-۱- مدل فضا برداری

مدل فضا برداری یک ابزار استاندارد در طول سه دهه اخیر برای سیستمهای بازیابی اطلاعات شده است [۱,۲,۵]. یکی از محاسن این روش امکان رتبه بندی مرتبط^۱ برای متون با پرس وجوی کاربران است. سیستم بازیابی اطلاعات فضا برداری [۳,۷] متون را به صورت بردارهایی از عبارات در نظر می گیرد. مجموعه متون شامل یک ماتریس $m \times n$ از m ترم و n متن است. ماتریس مذکور را ماتریس عبارت-متن می نامند که هر ستون نماینده یک متن از مجموعه و هر سطر نماینده یک عبارت است.

۱-۶-۱- مدل دودویی

در مدل برداری دودویی، هر متن به صورت بردار منطقی از تمام عبارات موجود در مجموعه است و مقدار هر ورودی یک است، اگر عبارت متناظر با آن ورودی در متن وجود داشته باشد و در غیر این صورت مقدار صفر خواهد داشت.

Terms	d1	d2	d2	q
↓	↓	↓	↓	↓
a arrived damaged delivery fire gold in of shipment silver truck	$A =$	$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}$	$q =$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$

شکل ۱: مدل برداری بولین از سه متن و یک پرس وجوی کاربر

در شکل ۱ مدل دودویی از سه بردار متن d_1 و d_2 و d_3 اندیس شده در ماتریس A ، و پرس وجوی q اندیس شده در بردار q ، را مشاهده می کنید. هر چهار بردار در فضای ۱۱ بعدی از عبارات اندیس شده اند. عیب این روش در این است که تعداد دفعات تکرار عبارات را در نظر نمی گیرد و فقط بود و نبود عبارت مهم است. معمولا برای رتبه بندی ارتباط پرس و جوی کاربر با متون از ضرب نقطه ای بین پرس و جو و هر یک از متون استفاده می شود.

^۱ Relevance Ranking

۲-۶-۱- مدل وزن دهی عبارت

در این مدل، متون را به صورت مخزنی از عبارات^۱ در نظر می گیرند. به این معنا که ترتیب عبارات در این روش اهمیتی ندارد. بر خلاف مدل بولیین، این مدل تعداد دفعات تکرار هر عبارت در متن را در نظر می گیرد و هرسلول ورودی از ماتریس، دفعات تکرار عبارت اندیس شده سطر متناظر سلول در متن اندیس شده ستون متناظر را نشان می دهد. یک متن d_i با بردار $x_i = \{f_1, f_2, \dots, f_m\}$ که یک فضای m بعدی از ترمها است نمایش داده می شود که در آن f_t دفعات تکرار عبارت t ام در متن i ام است. یکی از مهمترین محاسن این نوع مدل نمایش، ارائه ساختاری است که عملیات جبری را آسان و قابل درک می کند. وزن هر عبارت در هر متن تابعی است از دفعات تکرار عبارت در متن t (tf) و تعداد متونی از مجموعه (df) که عبارت در آن حداقل یک بار دیده شده است. در ادامه به تعدادی از مشهورترین این توابع اشاره می کنیم.

$$\begin{bmatrix} & D_1 & D_2 & \dots & D_n \\ T_1 & w_{11} & w_{12} & \dots & w_{1n} \\ T_2 & w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \dots & \dots & \dots & \dots \\ T_m & w_{m1} & w_{m2} & \dots & w_{mn} \end{bmatrix}$$

شکل ۲: مدل وزن دهی عبارت

در شکل ۲ نمایی کلی از مدل وزن دهی عبارت برای n متن و m ترم را نشان می دهد. وزن w معمولاً توسط یکی از روش های زیر محاسبه می شود.

وزن دهی $Tf-Idf$

تکرار عبارت ها (tf) و عکس تکرار متون (idf)، که اولی باعث افزایش اثر تکرار عبارت و دومی کاهش اثر عبارت هایی است که در اسناد زیادی تکرار شده اند. فرض بر این است عبارت هایی که در متون زیادی ظاهر شده اند خاصیت جداکنندگی بالایی ندارند و عمدتاً باعث کاهش صحت نتایج می شوند. این فرضیه در مورد مجموعه داده های متوازن که تعداد نمونه های همه خوشه ها تقریباً نزدیک به هم است خوب عمل می کند، اما برای مجموعه داده های نامتوازن ممکن است وزن برخی عبارت های مهم با قدرت جداکنندگی نسبتاً بالا را در خوشه های بزرگ به علت کوچک بودن عکس تکرار متون را کاهش دهد.

¹ Bag of Words

² Term Frequency

³ Document Frequency