



دانشگاه فردوسی مشهد
دانشکده مهندسی - گروه مهندسی کامپیوتر

پایان نامه کارشناسی ارشد

خوشه‌بندی اسناد متنی مبتنی بر مفاهیم همسایگی و شباهت معنایی

ملیحه دانش

استاد راهنما: جناب آقای دکتر محمود نقیب‌زاده

استاد مشاور: جناب آقای دکتر احد هراتی

تیر ۹۰

تقدیر و تشکر

خداوند را سپاس می‌گوییم که به من فرصت داد تا عمر خود را در راه تحصیل علم و دانش سپری کنم و همواره استادانی دلسوز و فرزانه بر سر راهم قرار داد تا در این راه دراز و بی‌پایان علم‌جویی، راهنمای راهم و تسکین آتش سیری‌ناپذیرم باشند. امید آن که به یاد خورشید تابان راهم، شمع کوچکی بر سر راه تشنگان دیگر باشم.

با تشکر از

راهنمای دلسوز و فرزانه

استاد ارجمند

جناب آقای دکتر محمود نقیب‌زاده

مشاور و مشوق راه علم

جناب آقای دکتر احد هراتی

و همه کسانی که تاکنون مرا در راه رسیدن به اهدافم یاری نموده‌اند.

چکیده

خوشه‌بندی، روش داده‌کاوی قدرتمندی است که جهت کشف موضوع از اسناد متنی مورد استفاده قرار می‌گیرد. در این زمینه الگوریتم‌های خانواده k-means به دلیل سادگی و سرعت بالا، در خوشه‌بندی داده‌هایی با ابعاد بالا، کاربرد فراوانی دارند. در این الگوریتم‌ها، معیار شباهت cosine، تنها شباهت میان زوج اسناد را اندازه‌گیری می‌کند که در مواقعی که خوشه‌ها به خوبی تفکیک نشده باشند، عملکرد مناسبی ندارد. درمقابل، مفاهیم همسایگی و اتصال با در نظر گرفتن اطلاعات سراسری در محاسبه میزان نزدیکی دو سند، عملکرد بسیار بهتری دارند. چنانچه میزان شباهت دو سند از حد آستانه‌ای بیشتر باشد آن دو سند همسایه‌اند و تعداد همسایه‌های مشترک میان آنها، مقدار تابع اتصال این دو سند را نشان می‌دهد. بنابراین با توجه به اینکه تنها دو حالت همسایگی و عدم همسایگی داریم که با صفر و یک نمایش داده می‌شوند، مقداری از اطلاعات را در مورد میزان شباهت میان اسناد از دست می‌دهیم که منجر به کاهش دقت خوشه‌بندی حاصل می‌شود. جهت رفع این مشکل، در گام اول لیستی از مقادیر گسسته را برای تعیین بازه‌ای از مقادیر آستانه به جای تنها یک مقدار، در نظر گرفتیم که به دنبال آن درجات متفاوتی از همسایگی، بر اساس میزان شباهت میان اسناد خواهیم داشت. همچنین جهت افزایش هر چه بیشتر دقت نتایج حاصل، از منطق فازی نیز بهره برده و مقدار شباهت میان اسناد را با استفاده از مقادیر عضویت فازی نمایش دادیم. به این ترتیب میزان همبستگی میان اسناد را با استفاده از منطق فازی بهبود داده و گام جدیدی در کاربردهای منطق فازی برداشتیم.

همچنین در این مدل، روابط معنایی میان کلمات نادیده گرفته شده و تنها اسنادی با واژگان مشابه با یکدیگر مرتبط شده‌اند. در این پروژه پایانی از آنتولوژی WordNet جهت ایجاد مدل جدید نمایش اسناد بهره بردیم، بدین صورت که در آن از روابط معنایی به منظور وزن‌گذاری مجدد بسامد کلمات در مدل فضای برداری اسناد استفاده شده است. سپس مفاهیم همسایگی و اتصال را بر روی مدل حاصل اعمال نمودیم. نتایج حاصل از اعمال روش‌های پیشنهادی و ترکیبات آنها بر روی مجموعه داده‌های متن واقعی، حاکی از عملکرد موثر و مناسب‌تر الگوریتم پیشنهادی ما نسبت به روش‌های پیشین می‌باشد و می‌تواند جایگزین خوبی برای الگوریتم‌های پیشین در امر خوشه‌بندی اسناد باشد.

فهرست مطالب

فصل ۱- ورود به مطلب.....	۱
۱-۱- مقدمه.....	۱
۲-۱- متن کاوی.....	۲
۱-۲-۱- تعاریف متن کاوی.....	۳
۳-۱- مراحل اصلی فرآیند متن کاوی.....	۴
۴-۱- کاربردهای متن کاوی.....	۵
۵-۱- خوشه بندی.....	۵
۶-۱- هدف از خوشه بندی.....	۶
۷-۱- کاربردهای خوشه بندی.....	۷
۸-۱- خوشه بندی در مقابل طبقه بندی.....	۸
۹-۱- رویه خوشه بندی.....	۹
۱-۹-۱- نمایش الگو.....	۹
۲-۹-۱- شباهت الگو.....	۹
۳-۹-۱- خوشه بندی یا گروه بندی.....	۱۰
۴-۹-۱- انتزاع داده.....	۱۰
۵-۹-۱- معیارهای آزمون خروجی الگوریتم.....	۱۰
۱۰-۱- چالش های الگوریتم های خوشه بندی.....	۱۱
۱۱-۱- خوشه بندی متن.....	۱۲
۱۲-۱- تعریف مساله.....	۱۲
فصل ۲- مبانی اولیه تحقیق.....	۱۵
۱-۲- مقدمه.....	۱۵
۲-۲- داده متنی.....	۱۵
۳-۲- پیش پردازش متن.....	۱۶
۱-۳-۲- جمع آوری داده های متنی.....	۱۷

۱۸ Collection Reader-۲-۳-۲
۱۸ Detagger-۳-۳-۲
۱۸Tokenization -۴-۳-۲
۱۹ حذف کلمات توقف.....۱-۴-۳-۲
۲۰ ریشه یابی۲-۴-۳-۲
۲۲ هرس کردن.....۵-۳-۲
۲۲ وزن گذاری کلمات۶-۳-۲
۲۲ روش های نمایش اسناد.....۴-۲
۲۳ مدل دودویی۱-۴-۲
۲۴ مدل برداری.....۲-۴-۲
۲۵ مدل احتمالاتی.....۳-۴-۲
۲۶ مقایسه روش های مدلسازی اطلاعات.....۴-۴-۲
۲۶ مدل فضای برداری در بازیابی اسناد.....۵-۲
۳۱ استفاده از دانش پس زمینه در خوشه بندی متون.....۶-۲
۳۱ خوشه بندی مبتنی بر آنتولوژی.....۱-۶-۲
۳۳ آنتولوژی۱-۱-۶-۲
۳۴ WordNet آنتولوژی۲-۱-۶-۲
۳۶ الگوریتم های خوشه بندی سند.....۷-۲
۳۶ روشهای سلسله مراتبی۱-۷-۲
۳۸ الگوریتم های افراز کننده.....۲-۷-۲
۳۹ مفهوم نزدیکترین همسایه ها.....۸-۲
۴۳ فصل ۳- مروری بر کارهای گذشته.....۳-۲
۴۳ ۱-۳- مقدمه.....۳-۲
۴۳ ۲-۳- استفاده از آنتولوژی در خوشه بندی اسناد.....۳-۲
۴۳ ۱-۲-۳- روش های وارد کردن آنتولوژی در نمایش متن.....۳-۲
۴۳ ۱-۱-۲-۳- افزودن مفاهیم.....۳-۲
۴۴ ۲-۱-۲-۳- جایگزین کردن کلمات با مفاهیم.....۳-۲
۴۴ ۳-۱-۲-۳- استفاده از بردار مفاهیم به تنهایی.....۳-۲
۴۴ ۲-۲-۳- مقایسه روش های استفاده از آنتولوژی در خوشه بندی اسناد.....۳-۲
۴۶ ۳-۳- معیارهای روابط معنایی۳-۲
۴۸ ۱-۳-۳- معیارهای مبتنی بر لبه۳-۲

۴۸ کوتاه ترین مسیر	۳-۱-۱-۱
۴۸ اتصالات وزن گذاری شده	۳-۱-۲-۱
۴۹ Wu and Palmer	۳-۱-۳-۳
۵۱ Hirst-St. Onge	۳-۱-۴-۱
۵۲ Li	۳-۱-۵-۱
۵۲ Chadorow و Leacock	۳-۱-۶-۱
۵۳ معیارهای مبتنی بر گره	۳-۲-۲
۵۴ Resnik	۳-۲-۱-۱
۵۴ Jiang-Conrath	۳-۲-۲-۲
۵۵ Lin	۳-۲-۳-۳
۵۵ معیارهای مبتنی بر ویژگی	۳-۳-۳
۵۵ Tversky	۳-۳-۱-۱
۵۶ پیشینه استفاده از همسایگی و اتصال در خوشه بندی	۳-۴-۱
۶۱ تعریف همسایگی و اتصال در خوشه بندی اسناد متنی	۳-۴-۱-۱
۶۲ انتخاب مراکز خوشه های اولیه بر اساس رتبه بندی	۳-۴-۲
۶۳ معیار شباهت بر مبنای توابع cosine و link	۳-۴-۳
۶۶ انتخاب یک خوشه جهت شکافتن بر اساس همسایه های مراکز	۳-۴-۴

فصل ۴- کارهای پیشنهادی..... ۶۸

۶۸ مقدمه	۴-۱-۱
۷۱ روش پیشنهادی	۴-۲-۱
۷۱ گام اول: استفاده از روابط معنایی در تعیین همسایگی میان اسناد	۴-۲-۱-۱
۷۱ مدل فضای برداری معنایی	۴-۲-۱-۱-۱
۷۲ محاسبه میزان وابستگی دو کلمه با استفاده از آنتولوژی	۴-۲-۱-۲
۷۶ گام دوم: افزایش دقت همسایگی	۴-۲-۲
۷۸ همبستگی با مقادیر گسسته	۴-۲-۲-۱
۷۹ همبستگی فازی	۴-۲-۲-۲
۸۲ متغیر فازی ورودی	۴-۲-۲-۱-۱
۸۳ متغیر فازی خروجی	۴-۲-۲-۲
۸۴ قوانین استدلال فازی در پروژه	۴-۲-۲-۳
۸۶ همبستگی فازی - گسسته	۴-۲-۳-۳

فصل ۵- پیاده سازی و نتایج..... ۸۷

۸۷ مقدمه	۵-۱-۱
----	-------------	-------

۸۷	۲-۵- روش پیاده سازی و محیط کار.....
۸۷	۱-۲-۵- مجموعه داده ها.....
۹۰	۲-۲-۵- روش ارزیابی.....
۹۱	۳-۵- نتایج خوشه بندی.....
۹۶	۱-۳-۵- نتایج حاصل از اعمال گام های همبستگی گسسته، فازی و فازی-گسسته بر الگوریتم KM:.....
۹۹	۲-۳-۵- نتایج حاصل از اعمال گام همسایگی معنایی در خوشه بندی با استفاده از الگوریتم KM:.....
۱۰۱	۳-۳-۵- نتایج حاصل از ترکیب گام های پیشنهادی در خوشه بندی با استفاده از الگوریتم KM:.....
۱۰۴	۴-۳-۵- نتایج حاصل از روش پیشنهادی و ترکیب گام های مختلف آن با استفاده از الگوریتم BKM:.....
۱۱۱	۴-۵- تحلیل زمانی روش پیشنهادی.....
۱۱۳	فصل ۶- نتیجه گیری و پیشنهادها برای کارهای آینده.....
۱۱۶	مراجع.....
۱۲۲	پیوست ها.....
۱۴۰	چکیده انگلیسی.....
۱۴۱	صفحه عنوان انگلیسی.....

فهرست جداول

۷۵ ۱-۳- ماتریس مدل برداری پایه برای اسناد
۷۶ ۲-۳- ماتریس مدل برداری معنایی جدید برای اسناد
۸۸ ۱-۵- مجموعه داده‌های Reuter
۸۹ ۲-۵- خلاصه‌ای از ویژگی‌های مجموعه داده
۹۲ ۳-۵- نتایج شباهت معنایی جفت کلمات در روش‌های مبتنی بر مسیر و مقایسه با روش پیشنهادی
۹۳ ۴-۵- شباهت معنایی جفت کلمات در روش‌های مبتنی بر نظریه اطلاعات و مقایسه با روش پیشنهادی
۹۷ ۵-۵- نتایج F-Measure حاصل از روش‌های همبستگی فازی، گسسته و فازی-گسسته با الگوریتم KM
۹۹ ۶-۵- نتایج F-Measure حاصل از روش‌های شباهت معنایی و همسایگی معنایی با الگوریتم KM
۱۰۲ ۷-۵- نتایج F-Measure حاصل از روش‌های پیشنهادی مختلف و ترکیبات آنها با الگوریتم KM
۱۰۵ ۸-۵- نتایج F-Measure حاصل از روش‌های همبستگی فازی، گسسته و فازی-گسسته با الگوریتم BKM
۱۰۵ ۹-۵- نتایج F-Measure حاصل از روش‌های شباهت معنایی و همسایگی معنایی با الگوریتم BKM
۱۰۶ ۱۰-۵- نتایج F-Measure حاصل از روش‌های پیشنهادی مختلف و ترکیبات آنها با الگوریتم BKM
۱۱۰ ۱۱-۵- نتایج precision حاصل از روش‌های پیشنهادی مختلف و مقایسه با روش‌های پیشین
۱۱۰ ۱۲-۵- نتایج recall حاصل از روش‌های پیشنهادی مختلف و مقایسه با روش‌های پیشین

فهرست اشکال

۳ شکل ۱-۱- نمونه‌ای از متن کاوی
۴ شکل ۲-۱- فرآیند متن کاوی
۶ شکل ۳-۱- نمونه‌ای خوشه‌بندی با استفاده از معیار فاصله به عنوان عدم شباهت بین داده‌ها
۸ شکل ۴-۱- تمایز طبقه‌بندی با خوشه‌بندی
۱۶ شکل ۱-۲- گام‌های پیش پردازش متن
۳۷ شکل ۲-۲- دندوگرام الگوریتم خوشه‌بندی سلسله مراتبی
۴۰ شکل ۳-۲- یافتن اسناد مشابه
۵۰ شکل ۱-۳- شبکه معنایی میان مفاهیم اولیه dog و cat
۵۷ شکل ۲-۳- انواع حدود آستانه برای خوشه‌بندی SNN
۵۹ شکل ۳-۳- سه پایگاه داده با چگالی‌های متفاوتی از نقاط
۶۰ شکل ۴-۳- قابلیت حصول از طریق چگالی
۷۳ شکل ۱-۴- یک نمونه آنتولوژی
۷۷ شکل ۲-۴- ماتریس شباهت و ماتریس همسایگی متناظر با آن
۸۰ شکل ۳-۴- توابع عضویت بر طبق متحنی‌های چندجمله‌ای
۸۳ شکل ۴-۴- مجموعه‌های فازی برای متغیر فازی ورودی
۸۳ شکل ۵-۴- مجموعه‌های فازی برای متغیر فازی خروجی
۸۵ شکل ۶-۴- نمونه‌ای از قوانین فازی به کار رفته در سیستم فازی پیشنهادی
۸۵ شکل ۷-۴- نمونه‌ای از خروجی به دست آمده توسط سیستم فازی طراحی شده
۹۴ شکل ۱-۵- تاثیر مقادیر متفاوت ضریب γ در مقدار ضریب همبستگی
۹۵ شکل ۲-۵- تاثیر مقادیر متفاوت ضریب α بر مقدار F-Measure در تابع همبستگی فازی-گسسته

VIII

- شکل ۳-۵- تاثیر روش‌های همبستگی گسسته، فازی و فازی-گسسته با استفاده از الگوریتم KM ۹۸
- شکل ۴-۵- تاثیر روش‌های شباهت معنایی و همسایگی معنایی با استفاده از الگوریتم KM ۱۰۰
- شکل ۵-۵- تاثیر روش‌های پیشنهادی مختلف و ترکیبات آنها با استفاده از الگوریتم KM ۱۰۳
- شکل ۶-۵- تاثیر روش‌های همبستگی گسسته، فازی و فازی-گسسته با استفاده از الگوریتم BKM ۱۰۷
- شکل ۷-۵- تاثیر روش‌های شباهت معنایی و همسایگی با استفاده از الگوریتم BKM ۱۰۸
- شکل ۸-۵- تاثیر روش‌های پیشنهادی مختلف و ترکیبات آنها با استفاده از الگوریتم BKM ۱۰۹
- شکل ۹-۵- میانگین زمان اجرای مورد نیاز روش‌های مختلف در یک بار اجرای حلقه الگوریتم KM ۱۱۲

فهرست اختصارات به کار رفته در متن

AHC	Agglomerative Hierarchical Clustering
BKM	Bisecting KMeans
DC	Discret Correlation
DL	Description Logic
FC	Fuzzy Correlation
KDT	Knowledge Discovery in Text
KM	KMeans
LCS	Least Common Subsummer
LSI	Latent Semantic Indexing
NLP	Natural Language Processing
SN	Semantic Neighbor
SVD	Singular Value Decomposition
SVSM	Semantic Vector Space Model
TF-IDF	Term Frequency-Inverse Document Frequency
VSM	Vector Space Model
WLC	Wordnet Lexical Categories

فصل ۱- ورود به مطلب

۱-۱- مقدمه

با ورود به دهه سوم جامعه جهانی وب، انقلاب متنی به عنوان یک تحول شگرف در قابلیت استفاده از اطلاعات موجود در اینترنت شده است. یافتن اطلاعاتی که تقریباً برای هر نیازی که پیش از این خودکار نبوده تنها با فشردن یک کلید و یا یک حرکت ماوس مقدور شده است. در مقایسه با انواع داده‌هایی که در پایگاه داده‌ها ذخیره می‌شوند، متن‌ها، داده‌هایی بدون ساختار، بی نظم و جهت بررسی الگوریتمی مشکل هستند. با این وجود در جامعه متمدن امروزی، متون به عنوان یکی از رایج‌ترین ابزارها جهت تبادل اطلاعات رسمی محسوب می‌شوند و داده‌های متنی بسیار پرکاربرد هستند. در حال حاضر بخش قابل توجهی از اطلاعات قابل دسترس در پایگاه داده‌های متنی (یا پایگاه داده‌های سند)^۱ که شامل مجموعه بزرگی از اسناد منابع مختلف (مثلاً مقالات خبری، مقالات، کتاب‌ها، ایمیل‌ها و صفحات وب) ذخیره شده‌اند. پایگاه داده‌های متنی به علت افزایش مقدار اطلاعات موجود به مدل الکترونیکی سریع رشد می‌کنند. امروزه بیشتر اطلاعات در صنعت، کسب و کار و سازمان‌های دیگر به صورت الکترونیکی و در قالب پایگاه داده متنی ذخیره شده‌اند.

به دنبال افزایش داده‌های متنی، مدیریت و تحلیل آنها از اهمیت بسیار زیادی برخوردار است. داده‌های ذخیره شده در بیشتر پایگاه داده‌های متنی، داده‌های نیمه ساخت‌یافته هستند، چون نه به طور کامل غیرساخت یافته و نه به طور کامل ساخت یافته‌اند [Sha 2005]. برای مثال یک سند شامل تعدادی حوزه ساخت یافته مانند عنوان، نویسندگان، تاریخ انتشار، رده^۲ و ... و از طرف دیگر شامل برخی اجزای متنی غیرساخت‌یافته مانند چکیده و محتویات است. روش‌های بازیابی اطلاعات^۳ (IR) مانند (شاخص‌گذاری متن) برای کنترل کردن سندهای غیر ساخت یافته ایجاد شده‌اند. روش‌های بازیابی اطلاعات قدیمی برای مقدار زیادی داده متنی که به طور فزاینده‌ای افزایش می‌یابند، ناکارآمد هستند. بدون دانستن محتویات سندها، فرموله کردن پرس و جوهای مناسب برای تحلیل و استخراج کردن اطلاعات مفید از داده، مشکل است. کاربرها نیاز به ابزارهایی برای مقایسه سندهای مختلف، مرتب کردن سندها بر اساس مرتبط بودن

¹ Document databases

² category

³ Information Retrieval

آن‌ها و یافتن الگوها دارند. بنابراین یکی از جدیدترین زمینه‌های مورد تحقیق در داده‌کاوی، "متن‌کاوی"^۱ برای این منظور گسترش یافت.

۱-۲- متن کاوی

متن‌کاوی فناوری ایجاد شده جهت کنترل داده‌های متنی در حال رشد است که در جهت برچینی اطلاعات معنی‌دار از متون زبان طبیعی تلاش می‌کند. متن‌کاوی یعنی جستجوی الگوها در متن غیرساخت یافته و برای کشف خودکار دانش مورد علاقه یا مفید از متن نیمه ساخت‌یافته استفاده می‌شود [Tan 2005].

متن‌کاوی تقریباً معادل با تجزیه و تحلیل متون است که وظیفه آن استخراج اطلاعات با کیفیت بالا از متن می‌باشد [Kan 2007]. در موارد معدودی نیز به عنوان فرآیند تحلیل متن جهت استخراج اطلاعاتی که برای اهداف خاصی مفید هستند، تعریف می‌شود. در زمینه کاوش متن معمولاً با متونی مواجه هستیم که وظیفه آنها ارتباط اطلاعات حقیقی و یا عقاید می‌باشد و هدف آن استخراج خودکار اطلاعات از چنین متونی می‌باشد، هر چند موفقیت جزئی حاصل گردد [Kan 2007].

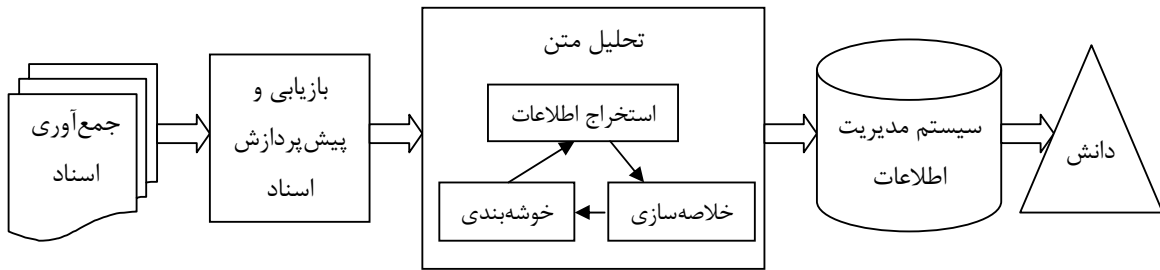
به طور کلی متن‌کاوی جهت مشخص کردن سیستمی که بتواند حجم زیادی از متون زبان طبیعی را تحلیل کند و الگوهای مفید زبانی و لغوی را شناسایی کرده و به دنبال آن اطلاعات احتمالاً مفید را استخراج کند، استفاده می‌شود [Fan 2005]. شکل ۱-۱ یک مدل کلی از یک کاربرد متن‌کاوی را نشان می‌دهد. این مدل با مجموعه‌ای از اسناد شروع می‌شود، یک ابزار متن‌کاوی، یک سند خاص را بازیابی و پیش‌پردازش می‌کند. سپس یک مرحله تحلیل متن انجام شده و در مواقعی از شیوه‌های مکرر تا استخراج اطلاعات استفاده می‌شود. سه روش تحلیل متن در این نمونه نشان داده شده اما بسیاری از ترکیبات دیگر نیز بر اساس اهداف سازماندهی می‌توانند استفاده شوند. اطلاعات حاصل می‌تواند در یک سیستم مدیریت اطلاعات قرار داده شود و در نهایت حجم وسیعی از دانش برای کاربر آن سیستم فراهم می‌شود [Fan 2005].

گاهی به جای واژه متن‌کاوی از عبارت "کشف دانش از متن" (KDT^۲)، استفاده می‌شود [Sha 2005]. معمولاً وظایف متن‌کاوی شامل طبقه‌بندی متن، خوشه‌بندی متن، استخراج مفهوم، تحلیل معنایی،

¹ Text Mining

² Knowledge Discovery in Text

خلاصه‌سازی متن و مدل‌سازی روابط میان نهادهای می‌باشد که در ادامه به طور خلاصه این مفاهیم را توضیح می‌دهیم.



شکل ۱-۱- نمونه‌ای از متن‌کاوی [Fan 2005]

۱-۲-۱- تعاریف متن‌کاوی

متن‌کاوی یا کشف دانش از متن، از شیوه‌های بازیابی اطلاعات، استخراج اطلاعات و پردازش زبان طبیعی (NLP^۱) استفاده کرده و آنها را به الگوریتم‌ها و روش‌های داده‌کاوی^۲، یادگیری ماشین^۳ و آماری^۴ مرتبط می‌کند [Mar 2009]. با توجه به ناحیه‌های تحقیق گوناگون، بر هر یک از آنها می‌توان تعاریف مختلفی از متن‌کاوی در نظر گرفت که دو نمونه از آن را در زیر می‌آوریم:

متن‌کاوی = استخراج اطلاعات: در این تعریف، متن‌کاوی متناظر با استخراج اطلاعات در نظر گرفته می‌شود (استخراج واقعیت‌ها^۵ از متن) [Wit 2000].

متن‌کاوی = کشف داده متنی: متن‌کاوی را می‌توان به عنوان روش‌ها و الگوریتم‌هایی از حوزه‌های یادگیری ماشین و آماری برای متن‌ها با هدف پیدا کردن الگوهای مفید در نظر گرفت. برای این هدف پیش‌پردازش کردن متون ضروری است. در بسیاری از روش‌ها، روش‌های استخراج اطلاعات، پردازش کردن زبان طبیعی یا برخی پیش‌پردازش‌های ساده برای استخراج داده از متون استفاده می‌شود. سپس می‌توان الگوریتم‌های داده‌کاوی را بر روی داده‌های استخراج شده اعمال کرد [Seb 2002].

¹ Natural language processing

² Data Mining

³ Machine Learning

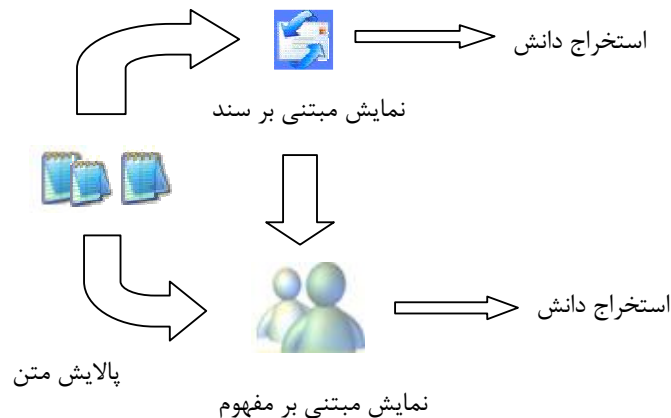
⁴ Statistic

⁵ Facts

در این پروژه پایانی بیشتر متن کاوی را به عنوان کشف داده متنی در نظر می‌گیریم و بیشتر بر روی روش‌های استخراج الگوهای مفید از متن برای دسته‌بندی مجموعه‌های متنی و یا استخراج اطلاعات مفید متمرکز می‌شویم.

۱-۳- مراحل اصلی فرآیند متن کاوی

متن کاوی فرآیندی است که شامل حوزه‌های تکنولوژیکی فراوانی است. بازیابی اطلاعات، داده کاوی، هوش مصنوعی و زبان‌شناسی محاسباتی، همگی حوزه‌هایی هستند که در این زمینه، نقش دارند. اما به طور کلی دو مرحله اصلی را در فرآیند متن کاوی می‌توان در نظر گرفت که در شکل ۱-۲ نشان داده شده است [Zhu 2004].



شکل ۱-۲- فرآیند متن کاوی [Zhu 2004]

اولین مرحله پیش پردازش مستندات است. خروجی این مرحله می‌تواند دو شکل مختلف داشته باشد: (۱) مبتنی بر سند^۱ (۲) مبتنی بر مفهوم^۲. در قالب نمایش مبتنی بر سند، آنچه که مهم است، نحوه‌ی نمایش بهتر برای مستندات است. مثلاً تبدیل اسناد به یک قالب میانی و نیمه ساخت یافته^۳، یا به کار بردن یک شاخص بر روی آن‌ها یا هر نوع نمایش دیگری که کار کردن با اسناد را کارآتر می‌کند. هر موجودیت^۴ در این نمایش در نهایت باز هم یک سند خواهد بود. در نوع دوم نمایش اسناد بهبود بخشیده می‌شود، مفاهیم و معانی موجود در سند و نیز ارتباط میان آنها و هر نوع اطلاعات مفهومی دیگری که قابل

¹ Document based

² Concept based

³ Semi-Structured

⁴ Entity

استخراج است، از متن استخراج می‌شود. در این نوع نمایش دیگر با مستندات به عنوان یک موجودیت مواجه نیستیم بلکه با مفاهیمی که از این مستندات استخراج شده‌اند، روبرو هستیم.

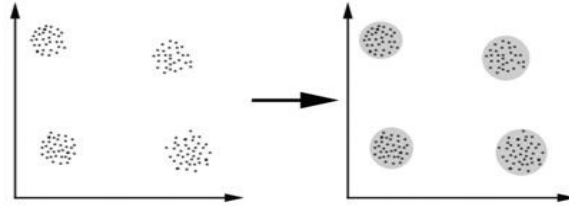
قدم بعدی استخراج دانش از این اشکال میانی نمایش اسناد است. بر اساس نحوه‌ی نمایش یک سند، روش استخراج دانش از یک سند متفاوت است. نمایش مبتنی بر سند، برای خوشه‌بندی، طبقه‌بندی، تجسم‌سازی و نظایر این‌ها استفاده می‌شود، در حالی که نمایش مبتنی بر مفهوم برای یافتن روابط میان مفاهیم، ساختن خودکار لغت‌نامه و آنتولوژی و نظایر آن بکار می‌رود [Zhu 2004].

۱-۴- کاربردهای متن‌کاوی

از جمله متداول‌ترین کاربردهای متن‌کاوی می‌توان موتورهای جستجو را نام برد که در آن کاربر یک عبارت یا کلمه (که ممکن است غلط املایی هم داشته باشد) را تایپ می‌کند و موتورهای جستجو توسط انباره بزرگی که از اسناد دارند، مرتبط‌ترین اسناد را پیدا می‌کنند [Zam 1997]. از جمله دیگر کاربردهای متن‌کاوی شامل طبقه‌بندی متن [Fab 2002]، خوشه‌بندی متن [Sal 1989]، استخراج مفهوم، تحلیل معنایی و خلاصه‌سازی متن [Raj 1997] را می‌توان نام برد. در این مستند، به کاربرد متن‌کاوی در خوشه‌بندی سند می‌پردازیم.

۱-۵- خوشه‌بندی

ما در جهانی پر از داده زندگی می‌کنیم. هر روزه انسان‌ها با حجم وسیعی از اطلاعات روبه‌رو هستند که باید آنها را ذخیره‌سازی کنند و یا نمایش دهند. یکی از روش‌های حیاتی کنترل و مدیریت داده‌ها، کلاس‌بندی یا گروه‌بندی داده‌های با خواص مشابه، درون مجموعه‌ای از دسته‌ها یا خوشه‌ها می‌باشد. خوشه‌بندی را می‌توان به عنوان مهم‌ترین مسئله در یادگیری بدون نظارت در نظر گرفت. خوشه‌بندی با یافتن یک ساختار درون یک مجموعه از داده‌های بدون برچسب درگیر است. خوشه به مجموعه‌ای از داده‌ها گفته می‌شود که به هم شباهت داشته باشند. در خوشه‌بندی سعی می‌شود تا داده‌ها به خوشه‌هایی تقسیم شوند که شباهت بین داده‌های درون هر خوشه حداکثر و شباهت میان داده‌های درون خوشه‌های متفاوت حداقل شود [Ros 2008]. نمونه‌ای از آن در شکل ۱-۳ آمده است.



شکل ۱-۳ - نمونه‌ای از اعمال خوشه‌بندی روی یک مجموعه از داده‌ها که از معیار فاصله به عنوان عدم

شباهت بین داده‌ها استفاده شده است [Ros 2008]

۱-۶- هدف از خوشه‌بندی

دسته‌بندی، جزئی از طبیعت انسان است. چیزهای محدودی وجود دارند که برای بقای زندگی بسیار مهم و ضروری هستند. تاریخ بشر به تاریخچه سرعت رشد دانش و پیاده‌سازی روش‌های تقسیم اشیاء به دسته‌های معلوم و قابل درک نیز برمی‌گردد. امروزه جهت بهبود توانمندی در دسته‌بندی، بشر بسیاری از روش‌های هوشمند را توسعه داده است. دو نمونه از موثرترین آنها شامل زبان‌های نوشتاری و کامپیوتر می‌باشد. با استفاده از هر دوی این امکانات ساختارهایی را توسعه داده و رده‌بندی‌های جدیدی را ایجاد کردیم.

در حال حاضر از کامپیوترها جهت توسعه و سازماندهی اطلاعات زیادی که یک نمونه از آنها به شکل متون هستند، استفاده می‌کنیم. ما انسان‌ها متون را به روش‌های مختلفی دسته‌بندی می‌کنیم. در اغلب موارد یک فرد ممکن است از یک مجموعه متون یکسان افزای بسیار متمایز و درعین حال مجاز و ارزشمندی را تولید کند. این امر ضرورتاً چیز بدی نیست. هر افزاز جدید از یک مجموعه از متون در صورتی که با استدلال درستی ایجاد شده باشد، ممکن است به ما دیدگاه و بصیرت جدیدی بدهد. افزازهای متون ممکن است در مواقعی منسوخ و بی‌ربط باشند. همچنین متون جدید ممکن است در یک ساختار قدیمی جای نگیرند یا اینکه ممکن است مجبور به ایجاد تغییراتی در یک ساختار شوند. ساخت یک افزاز جدید به طور دستی بسیار گران و زمانبر است. ابزارهای خودکاری که متون را افزاز می‌کنند، یا گروه‌های منطقی و معقولی را از متون استخراج کنند، می‌توانند بسیار ارزشمند باشند. حتی اگر افزازها و اجزای تولید شده از ابزارهای خودکار، نتایج بدتری از آنچه که به طور دستی و توسط بشر جمع‌آوری شده، داشته باشند باز هم ارزشمند هستند چرا که در اغلب موارد کسی تاکنون یک افزاز به طور دستی ایجاد نکرده است. از طرفی، اطلاعات مورد جستجو در محیط‌های الکترونیکی چنانچه توسط کاربرانی هدایت و بازیابی شوند، می‌توانند بهتر بازیابی شوند.

از یک دیدگاه، آنچه که به عنوان یک افراز خوب از یک مجموعه از متون محسوب می‌شود وابسته به استدلالی است که در ایجاد آن افراز استفاده شده که آیا صحیح است یا خیر و اینکه آیا به روش سازگاری استفاده شده یا خیر. در زمینه سازگاری، کامپیوتر بسیار ممتاز است اما استدلال باید به طریقی توسط انسان تامین شود. خوشه‌بندی با استفاده از روشی موثر در نمایش و تجسم اسناد و گروه‌بندی اسناد مشابه و مرتبط با هم و تشکیل مجموعه‌ای از اسناد، در این راستا گام موثری برداشته است. خوشه‌بندی روشی مفید جهت سازماندهی حجم وسیعی از اسناد متنی نامنظم در داخل تعداد محدودی از خوشه‌های معنی‌دار و منسجم می‌باشد که به موجب آن مکانیزم‌هایی جهت مرور و هدایت آگاهانه اطلاعات را تامین می‌کند.

۱-۷- کاربردهای خوشه‌بندی

خوشه‌بندی در زمینه‌های بسیاری کاربرد دارد، در زیر لیستی از مهم‌ترین آنها را یادآور می‌شویم:

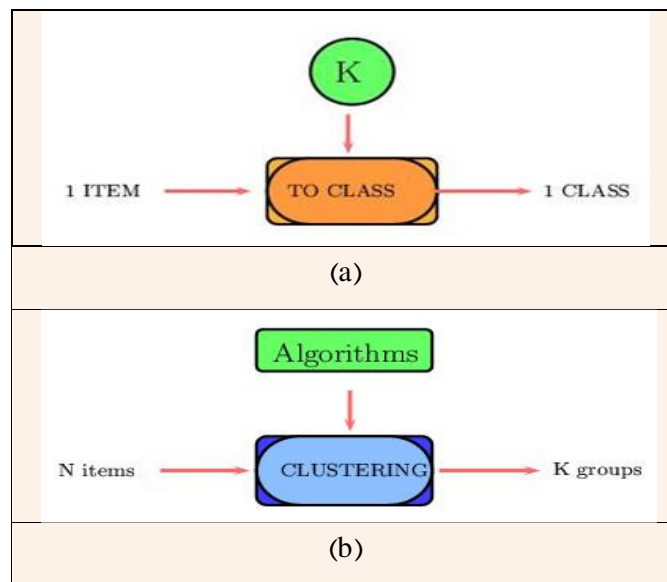
- در زمینه مهندسی (یادگیری ماشین، هوش مصنوعی، تشخیص الگو، مهندسی مکانیک و الکترونیک)
- علوم کامپیوتر (کاوش وب، تحلیل پایگاه داده فضایی، جمع آوری مستندات متنی، تقسیم بندی تصویر)
- علوم پزشکی (ژنتیک، زیست شناسی، میکروبی شناسی، فسیل شناسی، روان شناسی، بالین، آسیب شناسی)
- علوم زمین‌شناسی (جغرافیا، زمین‌شناسی، نقشه برداری از زمین)
- علوم اجتماعی (جامعه شناسی، روان‌شناسی، تاریخ، آموزش و پرورش)
- اقتصاد (بازاریابی، تجارت)

همچنین خوشه‌بندی ممکن است با نام‌های دیگری از قبیل علم رده‌بندی عددی [Gre 1996] و یادگیری بدون معلم (یا یادگیری بدون نظارت)^۱ [Wil 1988] بکار برده شود. امروزه، خوشه‌بندی نقش حیاتی در سازماندهی مجموعه‌های بزرگ مستندات متنی مانند وب دارد که در ادامه به آن می‌پردازیم.

^۱ Unsupervised Learning

۱-۸- خوشه‌بندی در مقابل طبقه‌بندی

با طبقه‌بندی خودکار، به ماشین این توانایی را می‌دهیم که تصمیم بگیرد یک متن، به کدام یک از مجموعه دسته‌های از پیش تعریف شده تعلق دارد. در خوشه‌بندی، ماشین تصمیم می‌گیرد که چطور یک مجموعه متن داده شده را باید تقسیم‌بندی کند. طبقه‌بندی زمانی مناسب است که شخص بخواهد متون جدیدی را بر طبق یک دسته‌بندی معلوم و مشخصی دسته‌بندی کند. خوشه‌بندی زمانی است که شخص بخواهد ساختارهای جدیدی را که پیش از این شناخته نشده بود، کشف کند. هر دو روش ممکن است که نتایج جالبی را بر روی داده‌های متنی نامعلوم ارائه بدهند، طبقه‌بندی آنها را مطابق با یک ساختار معلوم مرتب می‌کند و خوشه‌بندی ساختار این مجموعه معین را نمایش می‌دهد. در طبقه‌بندی هر داده به یک طبقه (کلاس) از پیشین مشخص شده تخصیص می‌یابد ولی در خوشه‌بندی هیچ اطلاعی از کلاس‌های موجود درون داده‌ها وجود ندارد و به عبارتی خود خوشه‌ها نیز از داده‌ها استخراج می‌شوند [Ros 2008].



شکل ۱-۴- (a) در طبقه‌بندی با استفاده یک سری اطلاعات اولیه داده‌ها به دسته‌های معلومی

نسبت داده می‌شوند. (b) در خوشه‌بندی داده‌ها با توجه به الگوریتم انتخاب شده به

خوشه‌هایی نسبت داده می‌شوند [Ros 2008].

بنابراین برخلاف طبقه‌بندی، در خوشه‌بندی گروه‌ها از قبل مشخص نمی‌باشند و همچنین معلوم نیست که برحسب کدام خصوصیات گروه‌بندی صورت می‌گیرد. در نتیجه پس از انجام خوشه‌بندی باید یک فرد خبره خوشه‌های ایجاد شده را تفسیر کند و در بعضی مواقع، لازم است که پس از بررسی خوشه‌ها بعضی از پارامترهایی که در خوشه‌بندی در نظر گرفته شده‌اند ولی بی‌ربط بوده یا اهمیت چندانی ندارند حذف