



پردیس علوم
دانشکده ریاضی، آمار و علوم کامپیوتر

شناسایی نوع سوال فارسی با استفاده از روش TF-IGM برای وزن دهی به عبارات

نگارنده

غزاله رستمی

استاد راهنما: دکتر هدیه ساجدی

پروژه برای دریافت درجه کارشناسی
در رشته علوم کامپیوتر

تیر ۱۳۹۷

چکیده

با توجه به حجم بالای داده های غیر ساختاری تولید شده توسط انسان که موجب عدم کارکرد مناسب موتور های جستجو در ارائه ی پاسخ مناسب به سوالات است نیاز به یک سیستم طبقه بندی سوال مطرح می شود. در یک سیستم طبقه بندی سوال هر سوال می تواند در یک ، چند و یا هیچ کلاسی قرار بگیرد . این موضوع می تواند در قالب یک یادگیری خودکار قرار گیرد تا بتوان به کمک آن هر سوال را به طور خودکار به دسته ای نسبت داد. هدف پروژه طراحی سیستمی برای طبقه بندی سوالات فارسی مبتنی بر روش وزن دهی TF-IGM به عنوان جایگزینی مناسب تر برای روش های سنتی TF-IDF است. این مدل یک مدل آماری جدید برای اندازه گیری قدرت تمایز میان دسته ها است. وزن دهی به عبارات مسئله ی مهمی در طبقه بندی متون است و به طور مستقیم روی دقت دسته بندی تاثیر دارد . از آنجا که روش TF-ID سنتی در دسته بندی چندان تاثیر گذار نیست روش های وزن دهی مختلفی توسط محققان پیشنهاد شده اند تاثیر گذاری این روش های پیشنهادی را با روش های دسته بندی مختلفی همچون ماشین بردار پشتیبان¹ و k نزدیک ترین همسایه² مورد بررسی قرار می دهیم . ابتدا عملیات پیش پردازش انجام می دهیم که شامل حذف کلمات کم معنی و بی معنی و علامت های نگارشی ، ریشه گیری ، بن واژه سازی و ... است . سپس با روش TF-IGM به ویژگی ها وزن تخصیص می دهیم قابل ذکر است که انتخاب ویژگی یک فاز در دسته بندی متون است و به روش های وزن دهی وابسته است ویژگی ها ی انتخاب شده باید در قالب مناسبی نمایش داده شوند در اینجا از مدل فضای برداری³ که یک مدل تاثیر گذار است بهره می بریم ، در این مدل هر مستند به صورت یک بردار از ویژگی ها که متناظر با آن برداری از وزن آنها داریم، نمایش داده می شود. در انتها به ارزیابی و ارائه ی نتایج می پردازیم.

¹Support Vector Machine (SVM)

²K-Nearest Neighbor (KNN)

³Vector Space Model (VSM)

پیشگفتار

پردازش دستی اطلاعات بدون ساختار به دلیل حجم زیاد و روز افزون آنها کاری زمان بر و پرهزینه است. متن کاوی عنوانی است که برای انجام این فرآیند خودکار انتخاب شده است. از جمله کاربردهای متن کاوی دسته بندی، خوشه بندی، خلاصه سازی و ... است. دسته بندی متون کاربردهای زیادی دارد در این پروژه به یکی از آن ها که دسته بندی سوالات فارسی است، اشاره می شود. طبقه بندی سوال به هیچ وجه کار بی اهمیتی نیست و نمی توان به سادگی و تنها با تکیه بر کلمات پرسشی به نتایج رضایت بخشی در زمینه ی دسته بندی سوالات رسید. ساختار های متنوع نحوی و جمله بندی سوالات دسته بندی آن ها را دشوار میکند. برای داشتن یک سیستم پاسخگویی کارا ابتدا به یک سیستم خوب برای دسته بندی سوالات نیاز داریم. در بخش دوم، پروژه به ارائه کارهایی می پردازد که تا کنون در زمینه ی دسته بندی سوالات فارسی و انگلیسی انجام شده است. در بخش سوم به ارائه روش پیشنهادی برای وزن دهی به عبارات و همچنین به پیاده سازی سیستم طبقه بندی سوال بر پایه ی روش پیشنهادی می پردازد. در بخش چهارم به مقایسه ی روش پیشنهادی با سایر روش های وزن دهی به عبارات می پردازد. در نهایت در بخش آخر به جمع بندی و نتیجه گیری حاصل از این پروژه پرداخته می شود.

فهرست مطالب

۱	مفاهیم مقدماتی	۱
۱	آشنایی با معماری موتور جستجو	۱.۱
۳	آماده‌سازی داده	۲.۱
۴	پاک‌سازی داده	۱.۲.۱
۴	معرفی روش‌های انتخاب ویژگی	۲.۲.۱
۵	استخراج ویژگی از سوال	۳.۲.۱
۷	آشنایی با برخی روش‌های دسته‌بندی	۳.۱
۸	کارهای پیشین در زمینه دسته‌بندی سؤال	۴.۱
۱۱	معرفی سیستم دسته‌بندی سوال بر مبنای روش وزن‌دهی پیشنهادی	۲
۱۱	معرفی مجموعه داده	۱.۲
۱۱	مدل فضای برداری	۲.۲
۱۳	کیف لغات	۱.۲.۲
۱۳	Word2Vec	۲.۲.۲
۱۵	معرفی روش پیشنهادی وزن‌دهی به ویژگی‌ها	۳.۲
۲۰	مقایسه‌ی روش پیشنهادی با سایر روش‌های وزن‌دهی به عبارات	۴.۲
۲۲	روش‌های ارزیابی و تفسیر مدل	۳
۲۲	معیارهای ارزیابی	۱.۳
۲۳	روش‌های ارزیابی	۲.۳
۲۳	روش Holdout	۱.۲.۳
۲۴	روش Random Subsampling	۲.۲.۳

۲۴ Cross Validation روش	۳.۲.۳
۲۴ Bootstrap روش	۴.۲.۳
۲۵	پیاده سازی و نتیجه گیری	۴
۲۵ آماده سازی و استخراج ویژگی از داده	۱.۴
۲۵ تولید بردار ویژگی	۲.۴
۲۶ TF-IDF روش	۳.۴
۲۷ استفاده از RapidMiner برای دسته بندی	۴.۴

فصل ۱

مفاهیم مقدماتی

۱.۱ آشنایی با معماری موتور جستجو

موتور جستجو نرم‌افزاری است که با گرفتن پرسشی^۱ از کاربر، مستنداتی را نمایش دهد که پاسخگوی نیاز اطلاعاتی وی باشد. پرسش به معنی کلمات کلیدی است که کاربر به کمک آنها نیاز اطلاعاتی خود را بیان می‌کند. یک موتور جستجو باید مستندات را در یک ساختار مناسب جمع‌آوری و نگهداری کند، بتواند مشابه‌ترین سند به پرسش کاربر را تشخیص دهد و مستندات را به ترتیب میزان ارتباط با پرسش کاربر نمایش دهد. در اینجا موتور جستجو به ابزارهایی اطلاق می‌شود که برای استخراج اطلاعات از صفحات وب، فایل متنی درون انباره، رکوردهای پایگاه داده استفاده می‌شوند. پردازش اولیه‌ای که موتور جستجو روی پرسش انجام می‌دهد شامل مراحل زیر است:

- بررسی عملگرهای موجود در پرسش
- حذف کلمات عمومی
- ریشه‌یابی کلمات
- استخراج کلمات کلیدی

¹Query

پس از این مراحل پرسش به خدمتگذار^۲ نمایه به منظور مقایسه با کلمات درون خدمتگذار ارسال می‌شود. بعد از انجام مقایسات، مستندات مرتبط شناسایی می‌شوند و مشخصات آنها به خدمتگذار اسناد ارسال شده تا به کاربر نمایش داده شود. موتور جستجو از سه بخش اساسی تشکیل شده است:

- گردآورنده اسناد
- نمایه ساز
- مدل های بازیابی اطلاعات والگوریتم دسته بندی

برخی از مشکلات زبان فارسی عبارتند از:

- گوناگونی معادل های علمی
- ضبط اسامی
- تعیین مرز کلمات (سرهم نویسی، جدانویسی، بی فاصله نویسی)
- انواع جمع ها
- صورت های مختلف نوشتاری
- استفاده از زبان محاوره در نوشتار

ریشه یابی^۳ یکی از پیچیده ترین مراحل کار در نمایه سازی متون است. اشتقاق یک واژه از ریشه اصلی موجب می‌شود تا برای ایفای نقش در جمله آماده شود. هدف ریشه یابی، زدودن الحاقات و یافتن جوهره اصلی واژه است. البته ممکن است در بعضی مواقع حذف الحاقات موجب از بین رفتن معنای واژه شود. از مزایایی که ریشه یابی دارد می‌توان به کاهش حجم ذخیره سازی نمایه و بهبود بازخوانی اسناد مرتبط با پرسش کاربر اشاره کرد. به عنوان مثال چند نمونه از سوال هایی که با نرم افزار پردازش زبان طبیعی دانشگاه فردوسی مشهد ریشه یابی کرده ایم را در جدول ۱.۱ مشاهده می‌کنیم

یکی از موارد مهم در نمایه ساز که نقش کلمات را از نظر میزان تاثیر آنها به عنوان کلمات کلیدی متن مشخص می‌کند وزن کلمه است. با استفاده از الگوهای مختلف وزن دهی، به هر واژه یا عبارت استخراج شده وزنی نسبت داده می‌شود که این وزن بیانگر میزان تاثیر کلمه در موضوع اصلی متن در مقایسه با سایر کلمات به کار رفته است.

²Server

³Stemming

سوال	سوال ریشه‌یابی شده
هیدراتهای کربن از کدام عناصر تشکیل شده اند	هیدراته کربن از کدام عنصر تشکیل کرد
به حجم پول، مجموع پول و شبه پول چه میگویند	به حجم پول، مجموع پول و شبه پول چه گفت
کوتاهترین جنگ در جهان چند روز طول کشید	کوتاه جنگ در جهان چند روز طول کشت

جدول ۱.۱: نمونه‌هایی از ریشه‌یابی سوالات فارسی

۲.۱ آماده‌سازی داده

آماده سازی داده به مجموعه عملیاتی که برای تولید داده‌های پالایش شده و قابل کاوش می‌شود آماده سازی داده گفته می‌شود. به بررسی دو نوع عملیات می‌پردازیم: [۶]

– استخراج داده

ابتدا داده در یک انباره داده، ذخیره می‌شود تا برای مراحل بعد مورد استفاده قرار بگیرد. این قسمت شامل دو مرحله است:

– جمع آوری داده

– انتخاب داده

– پیش‌پردازش داده

در این مرحله عملیاتی انجام می‌شود که موجب برطرف شدن مشکلات مختلف داده شده و برای انجام فرایند یادگیری آماده شود. از جمله این عملیات می‌توان به: پاک‌سازی داده، انتخاب زیر مجموعه ویژگی، فیلترینگ نمونه‌ها، نمونه‌برداری، کاهش ابعاد و... اشاره کرد. ما به معرفی دو نمونه از این عملیات می‌پردازیم:

۱.۲.۱ پاک سازی داده

یکی از مشکلات مهمی که داده می‌تواند داشته باشد این است که کیفیت آن کم باشد ابتدا مواردی که موجب کم شدن کیفیت داده می‌شود را معرفی می‌کنیم: نویز^۴، نمونه‌های پرت^۵، مقادیر از دست رفته^۶، داده‌های تکراری^۷.

نویز: مقدار اپسیلونی اگر به داده‌ی اصلی اضافه و یا از آن کم شود موجب می‌شود از داده‌ی اصلی دور بمانیم این تغییر یا تخریب را نویز می‌گوییم.

نمونه‌های پرت: نمونه‌هایی که مقادیر آنها با سایر مقادیر رکوردها بسیار متفاوت است. این نمونه‌ها کار الگوریتم‌های یادگیری را دچار مشکل می‌کنند.

مقادیر از دست رفته: ممکن است مقادیر برخی ویژگی‌ها null شود راه‌هایی برای مدیریت این مقادیر از دست‌رفته وجود دارد که می‌توان به حذف کردن، نادیده گرفتن، تخمین زدن و جایگزین کردن اشاره کرد.

داده‌های تکراری: رکورد‌هایی هستند که بار اطلاعاتی جدیدی ندارند و اطلاعات تکراری زیادی در آنها وجود دارد. حذف داده‌های تکراری می‌تواند اثرات مثبت و منفی در پی داشته باشد اثر مثبت به خاطر کاهش حجم داده و اثر منفی به این دلیل که می‌تواند موجب از بین رفتن نظم میان داده‌ها شود.

۲.۲.۱ معرفی روش‌های انتخاب ویژگی

در این بخش به طور مختصر به معرفی روش‌های انتخاب ویژگی می‌پردازیم. از جمله عملیات برای کاهش ابعاد می‌توان به انتخاب زیرمجموعه از ویژگی‌ها اشاره نمود. در این روش ویژگی‌های افزونه^۸ و غیر مرتبط حذف خواهند شد. ویژگی‌های افزونه ویژگی‌هایی هستند که با توجه به ویژگی‌های دیگر قابل محاسبه‌اند و موجب افزایش بی‌دلیل فضای الگوریتم می‌شوند. ویژگی‌های نامرتب و ویژگی‌هایی هستند که هیچ ارزش اطلاعاتی برای مسأله ندارند. [۶]

⁴Noise

⁵Outliers

⁶Missing Values

⁷Duplicate Data

⁸Redundant Features

روش ناآگاهانه : ^۹ در این روش همه زیرمجموعه‌های امکان‌پذیر از ویژگی‌ها به الگوریتم داده‌کاوی اعمال خواهند شد. به طور تصادفی یک نمونه کوچک از داده‌ها انتخاب می‌شود. این روش هنگامی که تعداد ویژگی‌ها کم باشد مناسب است. در واقع در این روش الگوریتم فقط یادگیری مدل را با تمام زیرمجموعه‌های ممکن از ویژگی‌ها را بر عهده می‌گیرد و کوششی برای یادگیری هوشمندانه انجام نمی‌دهد. [۶]

روش توکار : ^{۱۰} در این روش الگوریتم یادگیری مدل و انتخاب ویژگی را به صورت توأمان انجام می‌دهد. [۶]

روش فیلتری : ^{۱۱} در این روش قبل از اجرای الگوریتم، انتخاب ویژگی انجام می‌شود و الگوریتم فقط یادگیری مدل را انجام می‌دهد. [۶]

روش انحصاری : ^{۱۲} در این روش فقط الگوریتم وظیفه انتخاب ویژگی را دارد و یادگیری مدل را انجام نمی‌دهد. [۶]

یکی از کاربردهای مهم انتخاب زیرمجموعه ویژگی‌ها، تمرکز بر برخی ویژگی‌ها جهت یافتن ارتباط میان آنهاست. [۶]

۳.۲.۱ استخراج ویژگی از سوال

در حال حاضر برای دسته‌بندی سوالات ، انواع مختلفی از ویژگی‌ها وجود دارند که مورد استفاده قرار می‌گیرند از جمله‌ی این ویژگی‌ها که بر اساس اطلاعات زبان شناختی در نظر گرفته شده‌اند می‌توان سه دسته‌ی لغوی ^{۱۳} ، نحوی ^{۱۴} و معنایی ^{۱۵} را معرفی کرد .

ویژگی‌های لغوی معمولا کلماتی هستند که در سوال ظاهر می‌شوند. ویژگی‌های نحوی از ساختار نحوی یک سوال استخراج می‌شوند. دو نوع ویژگی نحوی رایج که در طبقه بندی سوال‌ها استفاده می‌شوند ، عبارتند از :

⁹Brute-Force Approach

¹⁰Embedded Approach

¹¹Filter Approach

¹²Wrapper Approach

¹³Lexical

¹⁴Syntactic

¹⁵Semantic

Tagged unigrams

به هر کلمه با توجه به نقش آن در سوال یک برچسب^{۱۶} اختصاص می‌دهیم. برای مثال در سوال ”رییس جمهور ایران در سال گذشته چه کسی بود؟“ داریم:
رییس جمهور-اسم، ایران-اسم، در-حرف اضافه، سال-اسم، گذشته-اسم، چه کسی-کلمه پرسشی، بود-فعل

Head words

یک سرکلمه به عنوان یک کلمه کلیدی در جمله در نظر گرفته می‌شود که حاوی اطلاعات مهمی برای شناسایی آن چه مورد جستجوی سوال است، می‌باشد. شناسایی درست سر کلمه می‌تواند دقت طبقه بندی را بهبود بخشد. برای مثال: ویژگی‌های معنایی در مورد داده‌های ناقص مفید هستند. با توجه به مفهوم معنایی در سطح بالا رابطه (یا شباهت معنایی) بین کلمات شناسایی می‌شوند. سه نوع ویژگی معنایی که در دسته بندی سوال استفاده می‌شود:

Question Category(QC)

برای اینکه تشخیص دهیم یک سوال به کدام دسته تعلق دارد باید از مقایسه شباهت سر کلمه آن سوال با همه ی دسته‌ها استفاده کنیم کلاس با بیشترین شباهت به عنوان یک ویژگی جدید در نظر گرفته می‌شود و به بردار ویژگی اضافه می‌شود مثلاً سوال کدام آهنگساز آمریکایی برای داستان غرب موسیقی نوشت؟ سر کلمه آهنگساز است بعد از مقایسه‌ی شباهت این کلمه با همه دسته‌های سوال متوجه خواهیم شد بیشترین نوع شباهت با دسته‌ی بشری است.

Question Expantion(QE)

یکی دیگر از ویژگی‌های معنایی گسترش سوال است. برای هر کلمه در سوال یک وزن تعریف می‌کنیم که با افزایش فاصله از سر کلمه مقدار وزن کاهش می‌یابد.

¹⁶Tag

Related Words(RW)

یک دیگر از ویژگی‌های معنایی کلمات مرتبط است . کلمات در گروه‌هایی قرار می‌گیرند که هر کدام با یک نام دسته نمایان می‌شوند . اگر یک کلمه در یک یا چند گروه وجود داشته باشد ، مقادیر مربوط آن به بردار ویژگی اضافه می‌شوند . مثلاً اگر هر یک از کلمات تولد ، تاریخ تولد ، روز ، دهه ، ساعت ، هفته ، ماه و سال در یک سوال باشند نام دسته (تاریخ) به بردار ویژگی اضافه می‌شود.

۳.۱ آشنایی با برخی روش های دسته بندی

ماشین بردار پشتیبان [۶]

استفاده از بردارهای پشتیبان خطی رویکرد جدیدی است که اخیراً مورد توجه بسیاری قرار گرفته است. این رویکرد دسته‌بندی با نظارت به این صورت عمل می‌کند که در مرحله‌ی آموزش سعی دارد مرز تصمیم‌گیری^{۱۷} را به گونه‌ای انتخاب نماید که حداقل فاصله‌ی آن با هر یک از دسته‌های مورد نظر را بیشینه کند. این امر باعث می‌شود که تصمیم‌گیری ما در عمل، شرایط نویزی را به خوبی تحمل کند و پاسخ‌دهی مناسب داشته‌باشد. الگوریتم‌های مبتنی بر ماشین بردار پشتیبان الگوریتم‌هایی هستند که سعی می‌کنند یک حاشیه^{۱۸} را بیشینه کنند.

الگوریتم، خطی را برای جداسازی کلاس‌های مثبت و منفی در نظر می‌گیرد که حاشیه کناری آن بیشتر باشد. لازم به ذکر است که ممکن است داده‌ها به گونه‌ای نباشند که بتوان با یک خط مستقیم آنها را دسته‌بندی کرد که در این شرایط برای جداسازی داده‌ها از منحنی استفاده می‌شود و با تبدیل غیرخطی، داده‌ها را به فضایی می‌بریم که بتوان آنها را با خط جدا نمود. ماشین‌های بردار پشتیبان دارای خواص زیر هستند:

- دسته بندی با حداکثر تعمیم
- رسیدن به بهینه سراسری تابع هزینه
- تعیین خودکار ساختار و توپولوژی بهینه برای دسته‌بند
- مدل کردن توابع تمایز غیرخطی با استفاده از هسته‌های غیرخطی

¹⁷Decision Boundary

¹⁸Margin

K نزدیکترین همسایه [۶]

این روش از جمله روش‌های مبتنی بر حافظه است که در دسته بندهای تاخیری^{۱۹} مورد استفاده قرار می‌گیرد. دسته‌بندهای تاخیری دسته بندهایی هستند که مرحله‌ی یادگیری مدل در آن‌ها به صورت مستقل وجود ندارد یعنی مدلی یاد گرفته نمی‌شود. در این دسته‌بند ها کل مجموعه رکوردهای آموزشی ذخیره می‌شوند وقتی یک رکورد جدید وارد می‌شود فاصله اقلیدسی آن از سایر رکوردها محاسبه می‌شود و سپس دسته رکوردی که نسبت به سایر رکوردها به رکورد جدید نزدیک‌تر است به آن تخصیص داده می‌شود. روش K نزدیک ترین همسایه یک گروه شامل K رکورد از مجموعه رکوردهای آموزشی که نزدیک‌ترین رکوردها به رکورد آزمایشی باشند را انتخاب کرده و بر اساس برجسب مربوط به آنها در مورد دسته رکورد تصمیم‌گیری می‌کند. استفاده از این الگوریتم نیازمند تعیین سه موضوع است:

- باید یک مجموعه رکورد داشته باشیم.

- یک معیار محاسبه شباهت داشته باشیم.

- مقدار K نیز مشخص باشد.

برای مسائل دسته بندی دودویی بهتر است K را عددی فرد در نظر بگیریم زیرا امکان پیروز شدن یکی از دو دسته را بالا می‌برد و برای مسائل چند دسته‌ای بهتر است K را بزرگ‌تر از تعداد دسته‌ها و متفاوت با عدد تعداد دسته‌ها از لحاظ زوج یا فرد بودن در نظر بگیریم.

۴.۱ کارهای پیشین در زمینه دسته بندی سؤال

برای سیستم طبقه بندی سؤال سه رویکرد پیشنهاد می‌شود :

۱. بر پایه قوانین اصلی

۲. یادگیری ماشین

۳. روش های ترکیبی

¹⁹Lazy Classifier

روش اول سعی می‌کند تا سوال‌ها را با قوانین دستی که از قبل ساخته شده‌اند متناظر کند تا بتواند نوع پاسخ سوال را شناسایی کند اما نوشتن این قوانین به صورت دستی کاری ملال‌آور است و باعث می‌شود سیستم نهایی بسیار خاص شود. یادگیری ماشین یک متد راحت‌تر به جای نوشتن قوانین برای دسته‌بندی سوال‌ها پیشنهاد می‌دهد که سیستم می‌تواند از روی داده‌های آموزشی یاد بگیرد و برای داده‌ی تست کلاس مناسب را پیش‌بینی کند این نوع سیستم‌ها می‌توانند خود را با شرایط جدید تطبیق دهند. سیستم‌های ترکیبی بسیار جدید هستند و استفاده از آن‌ها چندان رایج نیست. لی و روث^{۲۰} (۲۰۰۲) در یک مقاله‌ی علمی در زمینه‌ی پردازش طبیعی با استفاده از یک مجموعه‌ی ویژگی متنوع شامل ویژگی‌های نحوی و معنایی به عملکرد ۸۴.۲ درصد دست یافتند. [۴] کریشان^{۲۱} و همکاران (۲۰۰۵) از مفهوم خبر رسان استفاده می‌کنند که یک عبارت پیوسته (یک تا سه کلمه) داخل سوال است که برای طبقه‌بندی دقیق مورد استفاده قرار می‌گیرد. آن‌ها یک چارچوب فراشناختی اتخاذ کردند که در آن ابتدا یک مدل توالی برای طبقه‌بندی خبر رسان‌ها، آموزش دادند سپس ویژگی‌های خبر رسان پیش‌بینی شده را با ویژگی‌های کلی‌تر ترکیب کرده و به یک بردار بزرگ از ویژگی‌ها تبدیل کرده و به کمک ماشین بردار پشتیبان دسته‌بندی کردند و خبر رسان‌ها با استفاده از تجزیه سوال ورودی شناسایی شدند. این رویکرد به دقت ۸۶.۲ درصد برای دسته‌بندی سوال دست یافت در حالی که دقت دسته‌به‌کمک خبر رسان‌ها ۸۵ درصد بود. [۵]

هوانگ^{۲۲} و همکاران (۲۰۰۸) بر خلاف لی و روث که از مجموعه ویژگی‌ها بسیار غنی استفاده می‌کردند، پیشنهاد میکنند که از مجموعه ویژگی فشرده اما موثر استفاده شود. به طور خاص آن‌ها ویژگی سر کلمه^{۲۳} را پیشنهاد میکنند و در مدل‌های ماشین بردار پشتیبان خطی و بیشترین آنتروپی^{۲۴} به ترتیب به دقت‌های ۸۹.۲ و ۸۹.۰ درصد برای ۵۰ کلاس، دست یافتند. [۳]

²⁰Li and Roth

²¹Krishnan

²²Huang

²³Head Word

²⁴Maximum Entropy (ME)

اولالر ویلیامز^{۲۵} (۲۰۱۰) از مجموعه ویژگی‌های معنایی استفاده کرده است و نشان داده است که اطلاعات معنایی به تنهایی برای تولید یک سیستم دسته بندی سوالات با کارایی بالا کافی است و سیستم نهایی اش به دقت ۸۶.۶ درصد روی دیتا ست استاندارد UIUC میرسد.

²⁵Olalere Williams

فصل ۲

معرفی سیستم دسته بندی سوال بر مبنای روش وزن دهی پیشنهادی

در این قسمت مدلی جهت دسته بندی سوالات فارسی ارائه می شود سپس به ارزیابی کارایی آن پرداخته می شود.

۱.۲ معرفی مجموعه داده

در این گزارش از مجموعه داده (۲۰۱۶) UTQD که شامل ۱۱۷۵ سوال فارسی است، استفاده می کنیم . این سوالات به صورت دستی جمع آوری و برچسب گذاری شده اند و در قالب ۸ دسته برای طبقه بندی سوال مورد استفاده قرار می گیرند. تعداد سوالات در هر دسته در جدول ۱.۲ نشان داده شده است .

۲.۲ مدل فضای برداری

مدل فضای برداری یکی از مدل های بازیابی اطلاعات است که در سطح وسیعی به کار می رود. در این مدل، هر مقوله اطلاعاتی شامل متون ذخیره شده و هر تقاضای اطلاعاتی زبان طبیعی به صورت مجموعه بردارهایی از عبارات نگهداری می شوند. به طور نظری، این عبارات می توانند از واژگان کنترل شده انتخاب شوند. به خاطر وجود مشکلاتی در تهیه این واژگان، عبارات از متون استخراج می شوند. معمولا برای کاهش اندازه واژگان از ریشه واژه ها استفاده می شود. همچنین معمولا از واژه های بازدارنده نظیر the, of, an,, صرف نظر می گردد. از تمام واژه های موجود در مستندات ، یک مجموعه واژگان به وجود می آید.

تعداد سوالات	دسته
۷۰	مخفف
۱۲۹	موجودیت
۱۶۰	توصیف
۱۹۸	مکان
۲۵۹	بشری
۲۱۶	عدد
۳۶	لیست
۱۰۷	آیا

جدول ۱.۲: تعداد سوالات در هر دسته از مجموعه داده‌ی معرفی شده

هر مستند به صورت برداری از تمام واژگان نمایانده می‌شود. بعید است واژه‌هایی که فاقد بار معنایی هستند و به طور معمول در مستند یافت می‌شوند، اطلاعات مهمی ارائه دهند، بنابراین می‌توان این واژه‌ها را برای سرعت دادن به پردازش، حذف کرد. واژه‌های تکراری که می‌توان از آنها چشم پوشید فهرست واژه‌های غیرمجاز را می‌سازند. در حذف واژه‌های غیر مجاز، باید دقت زیاد به کار برده شود. برای مثال: چنانچه واژه‌های غیر مجاز در جمله: « to be or not to be » حذف شوند، این جمله غیر قابل بازیابی خواهد بود. مدل فضای برداری، شیوه‌ای است برای نمایش مستندات از طریق واژه‌های موجود در آنها. این مدل، یک تکنیک استاندارد در بازیابی اطلاعات است. بر اساس مدل فضای برداری، می‌توان تصمیم گرفت که کدام مستندات شبیه به یکدیگر و یا به کلید واژه‌های جستجو شبیه هستند. هم چنین برای بسیاری از روش‌های پردازش متن، نیاز به نمایش عددی کلمات و متون داریم تا بتوانیم از انواع روش‌های عددی حوزه یادگیری ماشین مانند اکثر الگوریتم‌های دسته بندی روی لغات و اسناد استفاده کنیم. در این جا نقش مهم این نحوه نمایش آشکار می‌شود. به طور کلی، می‌توان مزیت‌های اصلی مدل فضایی برداری را چنین بیان نمود:

۱. طرح وزن دهی به اصطلاح در این مدل، عملکرد بازیابی را بهبود می‌بخشد.
۲. استراتژی تطبیق جزئی این مدل، بازیابی مستندات را مجاز می‌شمارد که به شرایط جستجو نزدیک هستند.
۳. فرمول رتبه بندی کسینوسی آن، مستندات را بر طبق درجه تشابهی که به موضوع جستجو دارند، مرتب می‌کند.

۱.۲.۲ کیف لغات

فرض کنید فرهنگ لغتی داریم با N کلمه و لغت که به ترتیب الفبایی مرتب شده اند و هر لغت یک مکان مشخص در این فرهنگ لغت دارد. حال برای نمایش هر کلمه، برداری در نظر میگیریم با طول N که هر خانه آن، متناظر با یک لغت در فرهنگ لغت ماست که برای راحتی کار فرض می کنیم شماره آن خانه بردار، همان اندیس لغت مربوطه در این فرهنگ لغت خواهد بود. با این پیش فرض، برای هر لغت ما یک بردار به طول N داریم که همه خانه های آن بجز خانه متناظر با آن لغت صفر خواهد بود. در خود ستون متناظر با لغت عدد یک ذخیره خواهد شد. با این رهیافت، هر متن یا سند را هم می توان با یک بردار نشان داد که به ازای هر کلمه و لغتی که در آن به کار رفته است، ستون مربوط از این بردار برابر تعداد تکرار آن لغت خواهد بود و تمام ستون های دیگر که نمایانگر لغاتی از فرهنگ لغت هستند که در این متن به کار نرفته اند، برابر صفر خواهد بود.

به این روش نمایش متون، کیف لغات^۱ می گوئیم که بیانگر این است که برای هر لغت در کیف یا بردار ما، مکانی در نظر گرفته شده است. با این روش ما دو بردار عددی داریم که حال می توانیم از این دو در الگوریتم های عددی خود استفاده کنیم. با وجود سادگی این روش، اما معایب بزرگی بر آن وارد است. مثلاً اگر فرهنگ لغت ما صد هزار لغت داشته باشد، به ازای هر متن ما باید برداری صد هزار تایی ذخیره کنیم که هم نیاز به فضای ذخیره سازی زیادی خواهیم داشت و هم پیچیدگی الگوریتم ها و زمان اجرای آنها را بسیار بالا می برد. از طرف دیگر در این نحوه مدل سازی فقط کلمات و تکرار آنها برای ما مهم بوده است و ترتیب کلمات یا زمینه متن (اقتصادی، علمی، سیاسی و...) تاثیری در مدل ما نخواهد داشت.

۲.۲.۲ Word2Vec

روشی دیگر که توسط گوگل در سال ۲۰۱۳ پیشنهاد شده است و روشی بسیار کارآمد و مناسب برای نمایش لغات و متون و پردازش آنها است روش Word2Vec است که هدف از این بخش آشنایی اولیه با این روش قدرتمند، نمایش برداری کلمات است که می تواند در بسیاری از کاربردهای نوین پردازش متن مانند سنجش احساسات، جستجوی متون مشابه یا پیشنهاد اخبار یا کالای مشابه استفاده شود.

¹Bag of words

در این روش به کمک شبکه عصبی یک بردار با اندازه کوچک و ثابت برای نمایش تمام لغات و متون در نظر گرفته شده و با اعداد مناسب در فاز آموزش مدل^۲ برای هر لغت این بردار محاسبه می شود. در این بردار هر ستون، نمایشگر کلمه یا ویژگی خاصی نیست و فقط یک عدد را نمایش می دهد. برای افزایش دقت این روش، مجموعه داده اولیه که برای آموزش مدل مورد نیاز است، باید حدود چند میلیارد لغت را که درون چندین میلیون سند یا متن به کار رفته اند، دربرگیرد.

بعد از ایجاد بردارهای مرتبط با هر لغت، برای نمایش برداری هر متن یا خبر، می توان بردار تک تک کلمات به کار رفته در آن را یافته و میانگین اعداد هر ستون را به دست آورد که نتیجه آن یک بردار برای هر متن یا سند خواهد بود. سرعت این آموزش بسیار بالاست و در عرض چند ساعت و یا چند دقیقه (بسته به این که از کدام یک از دو الگوریتم آموزش آن استفاده کنیم) می توان حجم عظیمی از داده ها را به این الگوریتم داد و بردارهای لغات را ایجاد کرد. به طور مختصر، این الگوریتم برای ساخت بردارهای کلمات از یکی از دو روش Skip-gram و continuous bag-of-words (CBOW) استفاده می کند. این دو روش که هر دو یک شبکه عصبی ساده هستند که بدون وجود لایه پنهانی که در اغلب روشهای شبکه عصبی وجود دارد، به کمک چند قانون ساده، بردارهای مورد نیاز را تولید می کنند.

در روش کیف لغات پیوسته (CBOW)، ابتدا به ازای هر لغت یک بردار با طول مشخص و با اعداد تصادفی (بین صفر و یک) تولید می شود. سپس به ازای هر کلمه از یک سند یا متن، تعدادی مشخص از کلمات بعد و قبل آنرا به شبکه عصبی می دهیم (به غیر از خود لغت فعلی) و با عملیات ساده ریاضی، بردار لغت فعلی را تولید می کنیم (یا به عبارتی از روی کلمات قبل و بعد یک لغت، آنرا حدس می زنیم) که این اعداد با مقادیر قبلی بردار لغت جایگزین می شوند. زمانی که این کار بر روی تمام لغات در تمام متون انجام گیرد، بردارهای نهایی لغات همان بردارهای مطلوب ما هستند. روش Skip-gram برعکس این روش کار می کند به این صورت که بر اساس یک لغت داده شده، می خواهد چند لغت قبل و بعد آنرا تشخیص دهد و با تغییر مداوم اعداد بردارهای لغات، نهایتاً به یک وضعیت باثبات می رسد که همان بردارهای مورد بحث ماست.

²Training model

از لحاظ الگوریتمی این دو روش شبیه هم هستند با این تفاوت که CBOW لغات هدف را از روی لغات متن ورودی پیش‌بینی می‌کند ولی Skip-gram به صورت برعکس از روی لغات مرجوعه هدف، لغات ورودی را پیش‌بینی می‌کند. برعکس کردن این چرخه دلخواه به نظر می‌رسد ولی از لحاظ آماری CBOW تأثیر نرمی بر روی همه اطلاعات توزیعی دارد (با رفتاری شبیه به یک مشاهده بر روی کل متن) و در کل این روش می‌تواند روشی مفید برای استفاده در مجموعه دادگان کوچک‌تر باشد. اما Skip-gram با هر زوج محتوا-هدف به صورت یک مشاهده جدید رفتار می‌کند و در مجموعه دادگان بزرگ‌تر بهتر جواب می‌دهد.

۳.۲ معرفی روش پیشنهادی وزن دهی به ویژگی‌ها

وزن دهی به عبارت^۳ ها مسئله مهمی در دسته بندی متون و تأثیرگذار روی دقت دسته بندی است. از آنجا که روش TF-IDF^۴ سنتی روش چندان مؤثری نبوده است محققان روش های مختلفی برای وزن دهی به عبارت‌ها ارائه داده‌اند. در این بخش به مقایسه‌ی چند روش وزن دهی مختلف و معرفی روش TF-IGM^۵ می‌پردازیم. TF-IGM شامل یک مدل آماری است تا قدرت تمایز کلاس‌ها برای یک عبارت را به دقت اندازه بگیرد. در دسته بندی، مستندات در مدلی به نام مدل فضای برداری نمایش داده می‌شوند در این مدل هر مستند به صورت برداری عددی نمایش داده می‌شود که شامل وزن بسیاری از عباراتی که از متن استخراج کردیم، هستند. [۱]

بنابراین چگونگی وزن دهی به عبارات با یک روش مناسب مسئله‌ای اساسی در دسته بندی متون است و به طور مستقیم در دقت اثرگذار است. وزن هر عبارت بسیار وابسته به این است که یک عبارت چه میزان قدرت در تمایز دسته‌ها دارد. TF-IDF از این جهت که برچسب تعیین شده در مجموعه داده‌ی آموزش را مورد بررسی قرار نمی‌دهد، روش چندان مؤثری نیست. دیبل و سباستیان^۶ (۲۰۰۳) وزن دهی با نظارت^۷ را پیشنهاد دادند که دارای سه نسخه است: TF-CHI, TF-IG, TF-GR که به ترتیب با جایگزین کردن عامل IDF در TF-IDF با توابع انتخاب ویژگی information-gain, gain-ratio و χ^2 بدست می‌آیند.

³Term

⁴Term Fequancy-Inverse Document Fequancy

⁵Term Fequancy-Inverse Gravity moment

⁶Debole and Sebastiani

⁷Supervised Term weighting (STW)

اگرچه بعدها دریافتند که روش های وزن دهی با نظارت همواره برتر از TF-IDF نیستند، با این حال نتیجه بخش تر از آن به نظر می رسند. مسئله ی وزن دهی و مخصوصا وزن دهی با نظارت بسیار مورد توجه محققان قرار گرفت و شما های مختلفی از وزن دهی با نظارت توسط آن ها پیشنهاد شده است . شما های وزن دهی سنتی مانند TF و TF-IDF دودویی هستند اینکه به یک عبارت وزن بدهیم نتیجه بخش تر از این است که بر مبنای حضور یا غیابش در مستند مقدار های ۰ و ۱ قرار بدهیم چرا که بعضی از کلمات در مستند پرتکرار و بعضی کمیاب هستند و ما باید به کلمات پرتکرار وزن بیشتری بدهیم اما در این صورت و با روش TF ممکن است به عبارت هایی که پرتکرار هستند اما قدرت تفکیک بالایی ندارند، وزن زیادی بدهیم که برای رفع این نقص از فاکتور IDF استفاده می کنیم که در نهایت یعنی با روش TF-IDF به عبارتی با TF بالا و DF پایین وزن بیشتری تخصیص می دهیم.

فرمول محاسبه ی وزن عبارت t_k با روش TF-IDF به صورت زیر است (سباستیانی ۲۰۰۲):

$$w(t_k) = tf_k \cdot \log\left(\frac{N}{df_k}\right) \quad (۱.۲)$$

$$w(t_k) = tf_k \cdot \log\left(\frac{N}{df_k} + 1\right) \quad (۲.۲)$$

$$w(t_k) = tf_k \cdot \log\left(\frac{N}{df_k}\right) + 1 \quad (۳.۲)$$

که tf_k ، یعنی تعداد تکرار عبارت t_k در مستند و df_k ، یعنی تعداد مستنداتی که شامل عبارت t_k هستند . به خوبی می دانیم که دسته بندی متون یک روش یادگیری ماشین با نظارت است که به مجموعه ای از متون با دسته ی مشخص برای آموزش مدل یا دسته بندی نیاز دارد . اما TF-IDF در وزن دهی به عبارات از اطلاعات مربوط به برچسب مستندات داده ی آموزش استفاده نمی کند بنابراین وزن محاسبه شده یک عبارت به طور کامل بازتاب دهنده ی اهمیت آن در دسته بندی متون نیست . پس روش های وزن دهی با نظارت پیشنهاد شدند

که اغلب با جایگزینی فاکتور عمومی وزن‌دهی با یکی از تابع‌های انتخاب ویژگی بدست می‌آیند. بعداً محققان شما‌های وزن‌دهی با نظارت اصلاح شده را پیشنهاد کردند به عنوان مثال میتوان به روش TF-RF⁸ که توسط لن⁹ و همکارانش (۲۰۰۹) و هم چنین روش TF-Prob¹⁰ که توسط لیو¹¹ و همکارانش (۲۰۰۹) پیشنهاد شدند، اشاره کرد. روش TF-RF از توزیع مستندات و تعداد نمونه‌های کلاس مثبت و منفی در محاسبه‌ی وزن عبارات استفاده می‌کند. و فرمول آن به صورت زیر است:

$$w_{t_k} = tf_{t_k} \cdot \log\left(2 + \frac{a}{\max(1, b)}\right) \quad (4.2)$$

a تعداد مستندات در کلاس مثبت که شامل عبارت t_k هستند و b تعداد مستنداتی که در کلاس منفی شامل ترم t_k هستند را نشان می‌دهد. و فرمول روش وزن‌دهی TF-Prob هم به شکل زیر است:

$$w_{t_k} = tf_{t_k} \cdot \log\left(1 + \frac{a}{b} \cdot \frac{a}{c}\right) \quad (5.2)$$

c تعداد مستندات در کلاس مثبت که شامل عبارت t_k نیستند. دسته‌ی دیگری از شما‌های وزن‌دهی بر پایه‌ی ICF¹² نیز وجود دارند. برای مطالعه‌ی دقیق و جامع این شما‌ها می‌توانید به مقاله‌ی فلان رجوع کنید. TF-IGM روش وزن‌دهی استفاده شده در این پروژه است که در گروه شما‌های وزن‌دهی معنایی قرار می‌گیرد. در این روش، وزنی که به هر عبارت تخصیص می‌دهیم براساس شباهت معنایی آن عبارت با یک دسته خاص است و هر دسته با یک سری کلید واژه بیان می‌شود.

⁸Term Fequancy-Relevance Fequancy

⁹Lan

¹⁰Term Fequancy-Probablity

¹¹Liu

¹²Inverse Class Fequancy

IGM مدلی برای اندازه گیری قدرت تمایز دسته برای هر عبارت است . وزن هر عبارت با قدرت تمایز دسته آن عبارت بسیار مرتبط است برای تشخیص این قدرت تمایز به طور ابتدایی میتوان به نحوه ی توزیع آن عبارت در دسته ها توجه کرد .

به طور کلی هرچه توزیع یک عبارت در دسته ها یکنواخت تر باشد قدرت تمایز کمتری دارد . اما عبارتی که تنها در کلاس خاصی یا تعداد کمی از کلاس ها می آید قدرت تمایز بیشتری دارد و به چنین عباراتی باید وزن بیشتری تخصیص داد.

اگر یک عبارت تمرکز بیشتری روی توزیع بین کلاسی داشته باشد قدرت تفکیک بیشتری هم دارد . بنابراین میتوان یک عبارت را مطابق با تمرکز توزیع بین کلاسی یا غیریکنواختی آن وزن دهی کرد. از IGM برای اندازه گیری غیریکنواختی یا سطح تمرکز توزیع بین کلاسی یک عبارت، که بیان گر قدرت تفکیک دسته آن عبارت است ، استفاده میکنیم. برای محاسبه ی تمرکز توزیع بین کلاسی ترم t_k ابتدا فرکانس هایش در دسته های مختلف را به صورت نزولی مرتب می کنیم :

مرتب کردن است و m تعداد دسته هاست. هر چه تعداد دسته هایی که وقوع یک عبارت در آن ها متمرکز کمتر باشد مرکز گرانی^{۱۳} به سر سمت چپ نزدیک تر است. به خصوص زمانی که یک عبارت تنها در یک دسته اتفاق می افتد، مرکز گرانی در مکان $r=1$ است. و اگر توزیع بین دسته های یکنواخت باشد یعنی $f_{k1}=f_{k2}=\dots=f_{km}$ باشد مرکز گرانی در مکان $m/2$ قرار می گیرد. بنابراین موقعیت مکانی مرکز گرانی بیانگر تمرکز توزیع بین دسته های است. برای کل فرکانس وقوع یک عبارت در مجموعه داده هرچه تمرکز توزیع بین دسته ای بیشتر باشد فاصله ی مرکز گرانی از مبدأ کمتر خواهد شد. IGM برای اندازه گیری تمرکز توزیع بین کلاسی پیشنهاد شده است و به صورت زیر تعریف می شود :

$$igm(t_k) = \frac{f_{k1}}{\sum_{r=1}^m f_{kr} \cdot r} \quad (۶.۲)$$

¹³Gravity Center

فرمول (۴.۲) را میتوان به صورت زیر هم نوشت :

$$igm(t_k) = \frac{1}{\sum_{r=1}^m \frac{f_{kr}}{\max(f_{ki})} \cdot r} \quad (۷.۲)$$

و فرم نرمال شده ی IGM به صورت زیر است :

$$nigm(t_k) = \frac{(1+m) \cdot m \cdot igm(t_k) - 2}{(1+m) \cdot m - 2} \quad (۸.۲)$$

بازه ی مقدار IGM از $\frac{2}{((1+m) \cdot m)}$ تا ۱ است و زمانی که توزیع یکنواخت باشد حداقل است و وقتی که $f_{k1} > 0$ و $f_{k2} = \dots = f_{km} = 0$ باشد آنگاه حداکثر است. حال به جای فاکتور IDF سنتی در وزن دهی به عبارات یک فاکتور عمومی جدید بر پایه IGM به صورت زیر تعریف می شود :

$$w_g(t_k) = 1 + \lambda \cdot igm(t_k) \quad (۹.۲)$$

λ یک ضریب قابل تنظیم بین ۵.۰ تا ۹.۰ است. دو شمای وزن دهی به عبارت TF-IGM و RTF-IGM با فرمول های زیر تعریف می شوند :

$$w(t_k, d) = t f_{kd} \cdot \left(1 + \lambda \cdot \frac{f_{k1}}{\sum_{r=1}^m f_{kr} \cdot r} \right) \quad (۱۰.۲)$$

$$w(t_k, d) = \sqrt{t_{kd}} \cdot \left(1 + \lambda \cdot \frac{f_{k1}}{\sum_{r=1}^m f_{kr} \cdot r}\right) \quad (11.2)$$

۴.۲ مقایسه‌ی روش پیشنهادی با سایر روش های وزن دهی به عبارات

در مقاله‌ای^{۱۴} که کوین چن^{۱۵} و همکارانش راجع به روش وزن دهی TF-IGM نوشته اند، عملکرد هشت شمای وزن دهی مختلف در دسته بندی چند کلاسه با استفاده از دسته‌بند های SVM و kNN روی سه مجموعه داده‌ی 20 Newsgroups و Reuters-21578 و TanCorp با معیارهای ارزیابی Macro-averaged f1 و Micro-averaged f1 مورد بررسی قرار گرفته است و نتایج ارزیابی ها نشان می‌دهد که دو شمای وزن دهی -TF-IGM و RTF-IGM عملکرد بهتری نسبت به شماهای دیگر (TF-IDF, TF-CHI, TF-Prob, TF-RF) ، در دسته‌بندی چند کلاسه دارند .

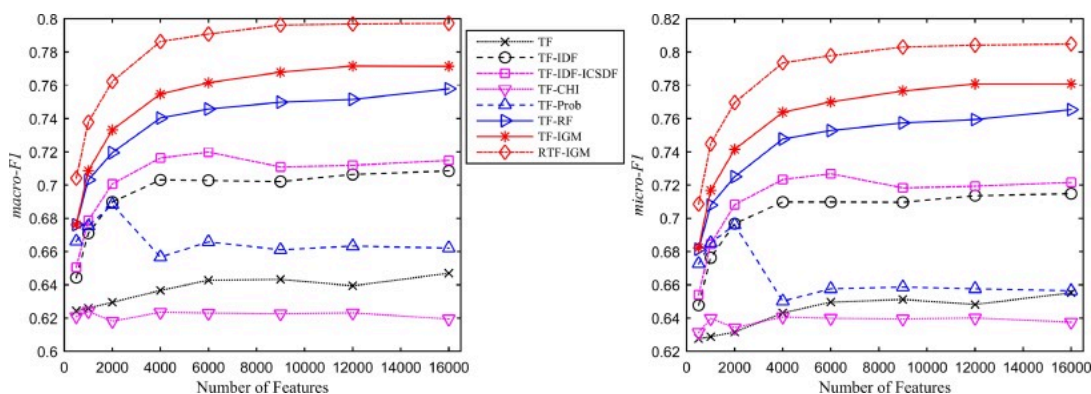
با این که TF-IDF معروف ترین شمای وزن دهی است و ثابت شده است که TF-RF کارایی بهتری نسبت به بسیاری از شماهای دیگر دارد شماهای پیشنهادی از هر دوی آن ها هم بهتر عمل میکنند.

به عنوان مثال در همین مقاله نتایج آزمایشات روی مجموعه داده‌ی 20 Newsgroups با دسته بند های SVM خطی و kNN(k=20) در شکل های ۱.۲ و ۲.۲ نشان داده شده است. هر منحنی در شکل ها بیانگر یک شمای وزن دهی است و محور عمودی نشان دهنده‌ی مقدار معیار های ارزیابی Macro-averaged f1 و Micro-averaged f1 در دسته‌بندی متون و محور افقی بیانگر تعداد واژه های متناظر است. واضح است که در دسته بندی با استفاده از SVM روی این مجموعه داده کارایی RTF-IGM و TF-IGM در هر تعداد از ویژگی ها از سایر شماهای نشان داده شده بهتر است . TF-RF و TF-IGSDF به مقدار بسیار کمی بهتر از TF-IDF هستند و TF عملکرد ضعیفتری نسبت به TF-IDF دارد و TF-CHI و TF-Prob عملکرد بسیار پایینی در مقایسه با بقیه دارند.

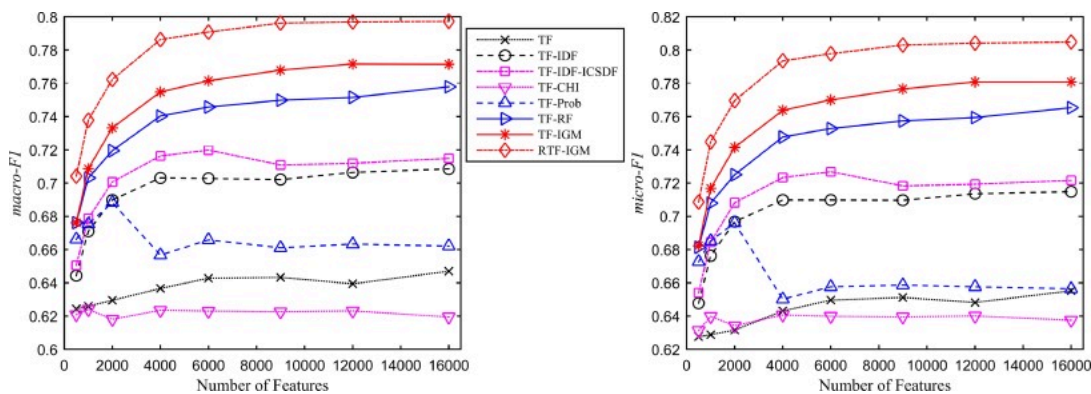
¹⁴Turning from TF-IDF to TF-IGM for term weighting in text classification

¹⁵Kewen Chen

در دسته بندی با استفاده از $kNN(k=20)$ همانطور که در شکل ۲.۲ نشان داده شده است تفاوت کارایی میان شماها بیشتر از تفاوتشان در حالت استفاده از دسته بند SVM است. برای هر تعداد ویژگی کارایی RTF-IGM و TF-IGM نسبت به دیگر شماها برتری دارد TF-RF و TF-ICSDF عملکرد بهتری از TF-IDF دارند. اما TF و TF-CHI و TF-Prob مقدار کمی ضعیف تر از TF-IDF عمل می کنند. [۱]



شکل ۱.۲: مقایسه‌ی کارایی هشت شمای وزندهی مختلف با دسته‌بند SVM روی مجموعه داده‌ی 20Newsgroups



شکل ۲.۲: مقایسه‌ی کارایی هشت شمای وزندهی مختلف با دسته‌بند $kNN(k=20)$ روی مجموعه داده‌ی 20Newsgroups

فصل ۳

روش های ارزیابی و تفسیر مدل

۱.۳ معیارهای ارزیابی

در آخرین مرحله پس از پیاده سازی الگوریتم های دسته بندی باید با استفاده از متون آزمایشی، صحت، دقت، بازخوانی، معیار ارزیابی F مدل پیشنهادی را بدست می آوریم. قبل از بیان روابط سنجش دقت دسته بندی نیاز به معرفی پیش نیازهای زیر است. ابتدا به جدول ۱.۳ توجه کنید :

- FP : تعدادی از داده ها که به غلط به عنوان دسته مثبت شناسایی می شوند.
 - TP : تعدادی از داده ها که به درست به عنوان دسته مثبت شناسایی می شوند.
 - FN : تعدادی از داده ها که به غلط به عنوان دسته منفی شناسایی می شوند.
 - TN : تعدادی از داده ها که به درست به عنوان دسته منفی شناسایی می شوند.
- معیار های ارزیابی دسته**

$$Recall = \frac{TP}{FN + TP} \quad (۱.۳)$$

$$Precision = \frac{TP}{FP + TP} \quad (۲.۳)$$

$$F_{measure} = \frac{2 \cdot Recall \cdot Precision}{Recall + Precision} \quad (۳.۳)$$

دسته مثبت		دسته منفی
FP	TP	دسته مثبت
TN	FN	دسته منفی

جدول ۱.۳: ماتریس درهم ریختگی برای مسئله‌ی دسته بندی دو دسته‌ای

$$micro-averagedF1 = \frac{2 * \sum_{i=1}^C |TP_{c_i}|}{2 * \sum_{i=1}^C |TP_{c_i}| + \sum_{i=1}^C |FN_{c_i}| + \sum_{i=1}^C |FP_{c_i}|} \quad (۴.۳)$$

$$macro - averagedF1 = \frac{1}{C} \sum_{i=1}^C F1_{c_i} \quad (۵.۳)$$

$$accuracy = \frac{TP + TN}{TN + TP + FN + FP} \quad (۶.۳)$$

$$Error = \frac{FP + FN}{TN + TP + FN + FP} \quad (۷.۳)$$

۲.۳ روش‌های ارزیابی

۱.۲.۳ روش Holdout

در این روش مجموعه داده به دو بخش با نام های داده آموزشی و آزمایشی تقسیم می‌شود. مدل دسته بندی توسط داده‌ی آموزشی ساخته شده و به وسیله ی داده آزمایشی ارزیابی می‌شود. چگونگی نسبت تقسیم به تشخیص تحلیل گر بستگی دارد و حسن این روش سادگی و سرعت بالاست. اولین ایراد این روش آن است که مجموعه داده‌ای که برای آزمایش استفاده می‌شود شانسی برای حضور در مرحله‌ی آزمایش ندارد. دومین ایراد این است که مدل بستگی به چگونگی تقسیم داده دارد اگر مجموعه داده‌ی آموزشی بزرگ در نظر گرفته شود دقت نهایی به دلیل کوچک شدن مجموعه داده‌ی آموزشی غیر قابل اعتماد خواهد بود.

۲.۲.۳ روش Random Subsampling

اگر روش Holdout را چندین بار اجرا کنیم و از نتایج حاصل میانگین گیری کنیم روش قابل اعتماد تری را برگزیده‌ایم مهمترین عیب این روش این است که در آن هیچ کنترلی بر روی تعداد دفعاتی که یک رکورد به عنوان نمونه‌ی آموزشی یا آزمایشی مورد استفاده قرار می‌گیرد وجود ندارد یعنی بعضی رکوردها ممکن است بیش از سایر رکوردها برای یادگیری یا ارزیابی به کار بروند.

۳.۲.۳ روش Cross Validation

در این روش کل مجموعه داده‌ها به k قسمت مساوی تقسیم می‌شوند. از $k-1$ قسمت به عنوان مجموعه داده‌های آموزشی استفاده می‌شود و براساس آن مدل ساخته می‌شود و با یک قسمت باقی مانده عملیات ارزیابی انجام می‌شود. فرآیند مزبور به تعداد k مرتبه تکرار خواهد شد، به گونه‌های که از هر کدام از k قسمت تنها یک‌بار برای ارزیابی استفاده شده و در هر مرتبه یک دقت برای مدل ساخته شده، محاسبه می‌شود. در این روش ارزیابی دقت نهایی دسته بند برابر با میانگین k دقت محاسبه شده خواهد بود.

۴.۲.۳ روش Bootstrap

در روش‌هایی که تا کنون گفته شد فرض بر آن است که انتخاب مجموعه‌ی آموزشی بدون جایگذاری صورت گرفته است. حال فرض می‌کنیم که هر رکورد مجدداً هم می‌تواند برای یادگیری مورد استفاده قرار گیرد. سپس رکوردهای انتخاب نشده برای ارزیابی مورد استفاده قرار می‌گیرند. این عملیات به تعداد b تکرار می‌شود احتمال انتخاب هر رکورد در مجموعه داده‌ی اولیه برابر با $1 - (1 - \frac{1}{N})^N$ است. اگر N به اندازه کافی بزرگ انتخاب شود، حاصل این رابطه 0.632 خواهد شد. به همین دلیل هر Bootstrap معادل 0.632 مجموعه داده‌ی اولیه خواهد شد.

فصل ۴

پیاده سازی و نتیجه گیری

۱.۴ آماده سازی و استخراج ویژگی از داده

در این بخش ابتدا عملیات پیش پردازش را روی مجموعه داده اعمال نمودیم هم‌چنین کلمات کم معنی و بی‌معنی^۱ را که در یک فایل متنی ذخیره کرده بودیم، از مجموعه داده حذف نمودیم. به تعداد کلاس‌ها فایل متنی ساختیم، نام هر فایل برچسب کلاس متناظر بود. سپس در هر فایل ویژگی‌های استخراج شده از تمام سوال‌های مجموعه داده که برچسب آن‌ها نام همان فایل بود، با در نظر گرفتن تکرار قرار داده شد. همه ویژگی‌های استخراج شده در فایلی به نام Features قرار گرفتند. از آنجایی که برای ساخت بردار داشتن ویژگی تکراری کار بی‌فایده‌ای است با حذف ویژگی‌های تکراری این فایل، فایل جدیدی تحت عنوان Selected features ساختیم.

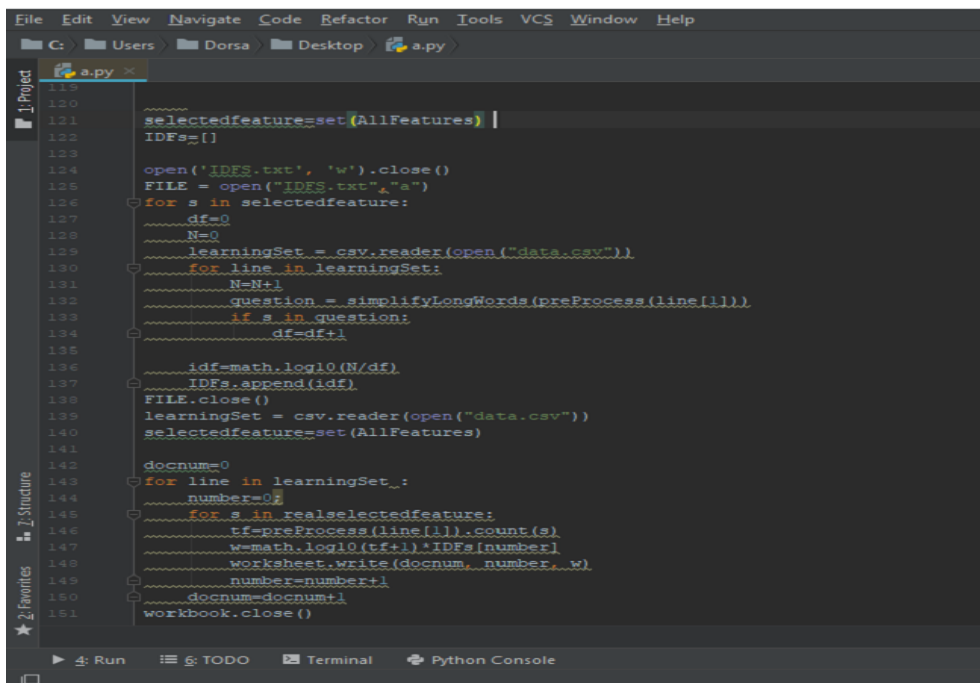
۲.۴ تولید بردار ویژگی

ماتریسی که قرار است برای دسته بندی از آن استفاده کنیم به تعداد سوالات مجموعه داده سطر و به تعداد ویژگی‌های انتخاب شده ستون دارد. اگر تنها به حضور یا عدم حضور ویژگی‌های انتخاب شده در هر سوال توجه می‌کردیم آنگاه اگر ویژگی انتخاب شده در سوال ظاهر شده بود عنصر متناظر در ماتریس را صفر و اگر ظاهر نشده بود، عنصر متناظر را یک قرار می‌دادیم.

¹Stop words

۳.۴ وزن‌دهی با روش TF-IDF

پر واضح است که این روش، روش مناسبی نیست زیرا تکرار ویژگی در سوال را در نظر نمی‌گیرد. ما در این‌جا از روش وزن‌دهی TF-IDF استفاده می‌کنیم تا این مشکل را برطرف کرده و دقت را بهبود بخشیم. همان‌طور در شکل ۱.۴ مشاهده می‌کنید ابتدا df محاسبه شده است سپس به ازای هر واژه از مجموعه داده آموزشی tf محاسبه شده است و طبق فرمول TF-IDF وزن هر واژه تولید شده است.



```
File Edit View Navigate Code Refactor Run Tools VCS Window Help
C:\Users\Dorsa\Desktop> a.py
119
120
121 selectedfeature=set(AllFeatures) |
122 IDFs=[]
123
124 open('IDFS.txt', 'w').close()
125 FILE = open("IDFS.txt","a")
126 for s in selectedfeature:
127     df=0
128     N=0
129     learningSet = csv.reader(open("data.csv"))
130     for line in learningSet:
131         N=N+1
132         question = simplifyLongWords(preProcess(line[1]))
133         if s in question:
134             df=df+1
135
136         idf=math.log10(N/df)
137         IDFs.append(idf)
138     FILE.close()
139     learningSet = csv.reader(open("data.csv"))
140     selectedfeature=set(AllFeatures)
141
142     docnum=0
143     for line in learningSet_:
144         number=0
145         for s in reselectedfeature:
146             tf=preProcess(line[1]).count(s)
147             w=math.log10(tf+1)*IDFs[number]
148             worksheet.write(docnum, number, w)
149             number=number+1
150         docnum=docnum+1
151     workbook.close()
```

شکل ۱.۴: کد پایتون برای وزن‌دهی به روش TF-IDF

در این مرحله مجموعه داده را با نسبت ۰.۷ برای داده آموزشی و ۰.۳ برای داده آزمایشی تقسیم می‌کنیم و پس از آن ماتریس آماده شده را به دسته بند های kNN و SVM می‌دهیم. کد پایتون مورد استفاده برای این دسته‌بندی‌ها را در شکل‌های ۲.۴ و ۳.۴ مشاهده می‌کنید. دقت با دسته‌بند kNN ، ۸۲.۲ درصد ، با دسته‌بند SVM ، ۷۹.۳ درصد است.

```

129
130
131
132
133
134 # KNeighborsClassifier
135 from sklearn.neighbors import KNeighborsClassifier
136
137 # Train Model
138 clf = KNeighborsClassifier(n_neighbors=15,algorithm='ball_tree').fit(X_train, y_train)
139
140 # Predict Model
141 y_predict_train = clf.predict(X_train)
142 y_predict_test = clf.predict(X_test)
143
144 # Compute Results
145 print('\nKNeighborsClassifier')
146 print('Accuracy on Train set: ',accuracy_score(y_train, y_predict_train))
147 print('Accuracy on Test set: ',accuracy_score(y_test, y_predict_test))
148
149
150
151

```

شکل ۲.۴: کد پایتون برای دسته‌بند kNN

۴.۴ استفاده از RapidMiner برای دسته‌بندی

در ادامه با نرم افزار رپیدماینر ماتریس در هم ریختگی را تولید کردیم که آن را در شکل ۴.۴ و ۵.۴ مشاهده می‌کنید. در فصل ۳ با معیارهای ارزیابی مدل آشنا شدیم. در این ماتریس معیارهای Precision و Recall برای هر کدام از دسته‌ها نمایش داده شده است. همان طور که مشاهده می‌شود این دو معیار برای دسته‌ی DESC نسبت به سایر دسته‌ها مقدار بیشتری را به خود اختصاص داده‌اند و این به این معناست که مدل ساخته شده روی این دسته عملکرد بهتری دارد. بالا بودن معیار Precision بیانگر این است که درصد بیشتری از سوال هایی که توسط مدل کلاس آن‌ها DESC پیش بینی شده واقعا متعلق به این کلاس هستند. هم‌چنین بالا بودن معیار Recall یعنی از میان سوال‌هایی که واقعا متعلق به کلاس DESC هستند، تعداد زیادی‌شان به درستی پیش بینی شده‌اند.


```

162
163
164
165
166
167
168 # Support Vector Machine
169 from sklearn import svm
170
171 # Train Model
172 clf = svm.SVC(kernel='rbf',cache_size=100).fit(X_train, y_train)
173
174 # Predict Model
175 y_predict_train = clf.predict(X_train)
176 y_predict_test = clf.predict(X_test)
177
178 # Compute Results
179 print('\nSupport Vector Machine')
180 print('Accuracy on Train set: ',accuracy_score(y_train, y_predict_train))
181 print('Accuracy on Test set: ',accuracy_score(y_test, y_predict_test))
182
183
184
185
186
187
188

```

شکل ۳.۴: کد پایتون برای دسته‌بند SVM

	true ENTY	true DESC	true ABBR	true HUM	true NUM	true LOC	class precision
pred. ENTY	335	4	2	10	65	111	63.57%
pred. DESC	1	291	0	17	0	2	93.57%
pred. ABBR	3	2	22	3	0	1	70.97%
pred. HUM	14	7	0	290	3	27	85.04%
pred. NUM	2	3	0	4	163	3	93.14%
pred. LOC	5	4	0	9	0	96	84.21%
class recall	93.06%	93.57%	91.67%	87.09%	70.56%	40.00%	

شکل ۴.۴: ماتریس در هم ریختگی داده آموزشی و دسته‌بند kNN

	true ENTY	true DESC	true ABBR	true HUM	true NUM	true LOC	class precision
pred. ENTY	53	2	0	4	0	6	81.54%
pred. DESC	0	96	0	5	0	1	94.12%
pred. ABBR	1	0	7	0	0	0	87.50%
pred. HUM	8	3	0	86	0	7	82.69%
pred. NUM	1	2	0	1	37	0	90.24%
pred. LOC	33	0	1	4	23	69	53.08%
class recall	55.21%	93.20%	87.50%	86.00%	61.67%	83.13%	

شکل ۵.۴: ماتریس در هم ریختگی داده آزمایشی و دسته‌بند kNN

منابع

- [1] Kewen Chen, Zuping Zhang, and Jun Long Hao Zhang. "Turning from TF-IDF to TF-IGM for term weighting in text classification." *Expert System With Applications*, vol.66 (2016) (Pages 245-260)
- [2] Mohammad Razzaghnoori, Hedieh Sajedi, and Iman Khani Jazani. "Question classification in Persian using word vectors and frequencies." *Cognitive systems research*, vol.47 (2018) (Pages 16-27)
- [3] Zhiheng Huang, Marcus Thint, and Zengchang Qin . "Question classification using head words and their hypernyms. "
- [4] Li, and D. Roth. "Learning Question Classifiers." *The 19th international conference on computational linguistics*, vol.1 (2002) (Pages 1-7)
- [5] V. Krishnan, S. Das, and S. Chakrabarti. "Enhanced Answer Type Inference from Questions using Sequential Models." *The conference on Human Language Technology and Empirical Methods in Natural Language Processing*.(2005)
- [6] صنیعی آباہه محمد، محمودی سینا، طاهرپرور محدثه. " داده کاوی کاربردی." تهران، انتشارات نیاز دانش(۱۳۹۴) (صفحات ۵۲-۵۶، ۱۲۳-۱۲۶)

Abstract

With the rapid growth of textual content on the Internet, automatic text categorization is a comparatively more effective solution in information organization and knowledge management. The necessity of the existence of Question Answering (QA) systems becomes evident by considering the fact that the enormous amount of unstructured data created by humans nowadays, results in ineffectiveness of search engines to provide the exact solution for a given question. Question classification plays an important role in question answering. Features are the key to obtain an accurate question classifier. Question classifier is a system that assigns a label to each question. Feature selection, one of the basic phases in statistical-based text categorization, crucially depends on the term weighting methods. In order to improve the performance of text categorization. Term weighting is a basic problem in text classification and directly affects the classification accuracy. Since the traditional TF-IDF (term frequency and inverse document frequency) is not fully effective for text classification, various alternatives have been proposed by researchers and here in this report we use the TF-IGM term weighting method in Persian question classification.



College of Science
School of Mathematics, Statistics, and Computer Science

Question classification in Persian using TF-IGM term weighting

Ghazale Rostami

Supervisor: Dr. Hedieh Sajedi

A thesis submitted to Graduate Studies Office
in partial fulfillment of the requirements for the degree of
B.Sc. in
Computer Science

2018