

چکیده

آنچه که بین تمام علوم مشترک است، حجم وسیعی از اطلاعات و داده به شکل‌های مختلف است، که همراه با پیشرفت و وسعت علوم بزرگتر و ارزشمندتر می‌شود. با رشد سریع علوم کامپیوتر و استفاده از آن در چند دهه اخیر، تقریباً همه سازمان‌ها در پروژه‌های مختلف خود حجم عظیمی داده در پایگاه-های داده ذخیره می‌کنند. این سازمان‌ها و کسانی که به نوعی در پروژه‌ها سهیم هستند به فهم این داده‌ها و بهتر بگوییم کشف دانش نهفته در آن نیازمندند. این نیاز، باعث به‌وجود آمدن حوزه جدید میان رشته‌ای کشف دانش و داده‌کاوی^۱ شده است، که حوزه‌های مختلف همچون پایگاه داده، آمار، یادگیری ماشین را برای کشف دانش ارزشمند نهفته در اطلاعات و داده‌ها را با هم تلفیق می‌کند. در این پژوهش به معرفی و بررسی مخازن نرم افزاری پرداخته شده است. که در فصل اول این پژوهش کلیات و تعاریف پژوهش قرار گرفته است. در فصل دوم ادبیات و مستندات پژوهش تشریح گردیده است و در فصل پایانی لزوم استفاده از تکنیکهای داده کاوی برای کشف و رفع خطاهای نرم افزار بر پایه ضریب تشابه معنایی متن و خوشه بندی بیان شده است.

^۱-Knowledge Discovery and Data Mining(KDD)

فهرست مطالب

فصل اول کلیات

۱-۱- مقدمه	۵
۲,۱. تعریف مسئله	۶
۳,۱. پیشینه تحقیقاتی	۷
۴,۱. فرضیه های تحقیق	۸
۵,۱. اهداف تحقیق	۸
۶,۱. نوآوری تحقیق	۹
۷,۱. ساختار پایان نامه	۹
۱,۷,۱. فصل اول: کلیات (همین فصل)	۹
۲,۷,۱. فصل دوم: ادبیات تحقیق	۹
۳,۷,۱. فصل سوم : استفاده از تکنیکهای دادهکاوی برای کشف و رفع خطاهای نرمافزار بر پایه ضریب تشابه معنایی متن و خوشه‌بندی	۱۰
۴,۷,۱. فصل چهارم: اعتبارسنجی و توسعه	۱۰
۵,۷,۱. فصل پنجم : نتیجه گیری	۱۰

فصل دوم ادبیات پژوهش

۱,۲. فرآیند مهندسی نرمافزار	۱۲
۱,۲,۱. پیاده سازی، آزمایش، تست و مستند سازی	۱۲
۲,۲. انواع مخازن داده	۱۳
۱,۲,۲. کد اصلی:	۱۳
۲,۲,۲. مخازن خطا (سیستم ردیابی خطا BTS)	۱۳
۴,۲,۲. پایگاه داده کنترل منبع (پایگاه داده کنترل ویرایش ها)	۱۴
۵,۲,۲. اطلاعات طراحی و نیازمندیهای سیستم	۱۵
۶,۲,۲. مخازن شرح اجرا	۱۵
۷,۲,۲. مخازن سیستم نرم افزار	۱۵
۳,۲. خطای نرمافزاری	۱۶

۱,۳,۲	سیستم ردیابی خطا	۱۶
۴,۲	تحقیقات پیشین در حوزه داده‌کاوی در مخازن خطا	۲۳
۵,۲	اندازه گیری شباهت بین متون	۲۶
۱,۵,۲	شباهت خطی	۲۶
۱,۱,۵,۲	اندازه گیری شباهت بر پایه کاراکتر	۲۷
۲,۵,۲	تشابه بر پایه مجموعه	۳۰
۳,۵,۲	تشابه بر پایه دانش	۳۳
۴,۵,۲	اندازه گیری شباهت ترکیبی	۳۴

فصل سوم استفاده از تکنیکهای داده کاوی برای کشف و رفع خطاهای نرم افزار بر پایه ضریب تشابه معنایی متن و خوشه بندی

۱,۳	مقدمه	۳۸
۲,۳	محاسبه تشابه معنایی بر پایه مجموعه و تشابه خطی	۳۹
۱,۲,۳	شباهت خطی بین کلمات	۴۰
۲,۲,۳	تشابه معنایی بین کلمات	۴۳
۳,۲,۳	تشابه جملات بر اساس عبارات مشترک	۴۹
۴,۲,۳	شباهت کلی جملات	۵۰
۵,۲,۳	نمونه توضیح روند کار	۵۲
۳,۳	خوشه بندی داده ها	۵۴
۱,۵,۲	الگوریتم K-MEANS	۵۴
	منابع	۵۶

کلیات

آنچه که بین تمام علوم مشترک است، حجم وسیعی از اطلاعات و داده به شکل‌های مختلف است، که همراه با پیشرفت و وسعت علوم بزرگتر و ارزشمندتر می‌شود. با رشد سریع علوم کامپیوتر و استفاده از آن در چند دهه اخیر، تقریباً همه سازمان‌ها در پروژه‌های مختلف خود حجم عظیمی داده در پایگاه‌های داده ذخیره می‌کنند. این سازمان‌ها و کسانی که به نوعی در پروژه‌ها سهیم هستند به فهم این داده‌ها و بهتر بگوییم کشف دانش نهفته در آن نیازمندند. این نیاز، باعث به وجود آمدن حوزه جدید میان رشته‌ای کشف دانش و داده‌کاوی^۱ شده است، که حوزه‌های مختلف همچون پایگاه داده، آمار، یادگیری ماشین را برای کشف دانش ارزشمند نهفته در اطلاعات و داده‌ها را با هم تلفیق می‌کند. اصطلاح Mining Data همان‌طور که از ترجمه آن به معنی داده‌کاوی مشخص می‌شود، به مفهوم استخراج اطلاعات نهان، و یا الگوها و روابط مشخص در حجم زیادی از داده‌ها در یک یا چند بانک اطلاعاتی بزرگ است.

هر نرم‌افزار در طول فرآیند تولید و پس از آن، انبوهی از اطلاعات و مستندات دارد که قابل کاوش و استفاده مفید است. این داده‌ها معمولاً در پایگاه داده‌هایی به نام مخازن نرم‌افزاری^۲ ذخیره و نگهداری می‌شوند. مخازن نرم‌افزاری نمایش دقیقی از مسیرتولید یک سیستم نرم‌افزاری ارائه می‌دهند [۱]. هدف از کاوش مخازن نرم‌افزاری MSR^۳ استفاده هوشمند از تحلیل داده‌های نهفته در آن‌ها برای کمک به تصمیم‌گیری‌های بهتر و سریع‌تر در پروژه تولید و پشتیبانی آن‌ها است. آنچه که در اینجا مورد توجه این تحقیق است، استخراج اطلاعات مهم برای همه ذینفعان پروژه نرم‌افزاری است. این اطلاعات از مجموعه داده‌های مرتبط با خطاهای رخ داده در طول تولید و پشتیبانی پروژه استخراج می‌شود.

در سال‌های گذشته مدل‌های مختلفی با استفاده از الگوریتم‌های داده‌کاوی، تشابه متن و دسته‌بندی و خوشه‌بندی داده‌ها ارائه شده. اما از آن‌جا که جستجو و استخراج اطلاعات از میان داده‌های متنی نیازمند روشی هوشمند برای تطبیق جنبه‌های معنایی و دستوری است، نیاز به مدل‌هایی که از الگوریتم‌ها معنایی استفاده کنند وجود دارد. در تحقیق‌های مورد مطالعه این نیاز حس می‌شود.

^۱-Knowledge Discovery and Data Mining(KDD)

^۲-Software Repositories

^۳-Mining Software Repositories

سعی بر این شده که با استفاده از الگوریتم معنایی برپایه مجموعه تشابه جملات^۱ بر پایه LCS^۲ [۳] و تشابه کلمات (SOC-PMI)^۳ [۳]، روی مستندات ذخیره شده در مخازن خطای نرم افزار، مانند راه حل های ارائه شده برای خطاهای مشابه مدل های قبلی را تکمیل کرده و جوابی بهینه و سریع تر برای خطای پیش آمده پیدا کنیم. همچنین می توان زمانی تخمینی نیز برای تصحیح خطا پیش بینی کرد تا راهنمای تیم توسعه و ذینفعان دیگر نرم افزار باشد. همچنین دید بهتری نسبت به روند پیشرفت و تکامل نرم افزار مورد نظر ارائه شود.

۲.۱. تعریف مسئله

یکی از مراحل مهم و اساسی در مهندسی و تولید نرم افزار مرحله یافتن و رفع خطاهای موجود در نرم افزار است. این مرحله از تولید نرم افزار جزء وقت گیرترین و پرهزینه ترین مراحل به حساب می آید [۴]. سال هاست که دانش داده کاوی و استخراج دانش به کمک مهندسين نرم افزار آمده است. رفع خطا در فرآیند تولید بسته به مدل توسعه نرم افزار چندین بار انجام می گیرد. خطاها و مشکلات برطرف شده معمولاً به روش های مختلف تحت عنوان مخازن خطای نرم افزار، مستندسازی و ذخیره می شود. این مخازن منابع عظیم دانش هستند، که کمک بزرگی در تسریع زمان تولید نرم افزار و پایین آوردن هزینه ها خواهد بود [۵]. روش هایی نیاز است که این دانش و اطلاعات مفید استخراج شود. در این تحقیق روشی برای سرعت بخشیدن به رفع خطای جدید با استفاده از اطلاعات موجود در مخازن خطای نرم افزار، ارائه شده است. مدل های زیادی تا به حال ارائه شده که یا مکمل هم بوده یا از الگوریتم های جدید استفاده شده است. مدل های پیشنهادی با استفاده از تشابه متن همگی از الگوریتم های معمولی و ساده استفاده کرده اند. در این مدل ها به این نکته مهم کمتر توجه شده که مخازن حجم بالا و پیچیده ای از اطلاعات را شامل می شوند، که بعضاً تشابه بین کلمات و معانی مختلف یک جمله نتیجه گیری را سخت تر می کند. پس نیاز به الگوریتم های معنایی در بررسی تشابه متن احساس می شود. همچنین باید به این نکته مهم توجه کرد که الگوریتم معنایی انتخاب شده بهینه است و قادر باشد میان این حجم اطلاعات که از سوی کاربران مختلف ثبت می شود، بهترین جواب با تشابه بیشتر را انتخاب کند. اگرچه تا به حال روش های بسیاری برای تکمیل مخازن خطا و استفاده از دانش نهفته در آن صورت گرفته اما می توان گفت که ضرورت اهمیت به تشابه معنایی بین داده ها در

^۱-String Similarity

^۲-longest common subsequence

^۳-Second-order co-occurrence pointwise mutual information

نظر گرفته نشده است. در این تحقیق سعی شده که این ضعف در جستجو و بهره‌گیری دانش نهفته در این مخازن داده پوشش داده شود.

مدل ارائه شده در اینجا ابتدا لیستی از خطاهای مشابه خطای جدید با استفاده از یک الگوریتم تشابه معنایی مناسب، با توجه به اطلاعات متنی ذخیره‌شده ارائه می‌دهد. در مرحله بعد این خطاها براساس چرخه عمر خطا با استفاده از روش خوشه‌بندی K-means، خوشه‌بندی می‌شوند.

همچنین میانگین تشابه هر دسته به خطای جدید، گروه منتخب را مشخص می‌کند. خطاهای موجود در این گروه، راه‌حل‌های پیشنهادی برای هر کدام، کمکی برای تسهیل و تسریع در رفع خطا است و میانگین طول عمر گروه تخمینی بر پیچیدگی و زمان حل مشکل خواهد بود.

۳.۱. پیشینه تحقیقاتی

شاید بتوان گفت که بنیان MSR و نرم‌افزارهای ردیابی خطا^۱ با تولید GNATS^۲ در سال ۱۹۹۲ نهاده شد. پس از آن نرم‌افزارهای متعددی به یاری مهندسی نرم‌افزار شتافته‌است. در این مدت سعی و اهتمام محققان بر تکمیل مدل‌ها و موتورهای جستجو مورد استفاده در این نرم‌افزارها بوده است. در سال ۲۰۰۰ مدلی با استفاده از رگرسیون فازی^۳ برای تخمین تعداد خطاهای موجود در نرم‌افزار ارائه شد [۶]. در [۵] همین کار با استفاده از شبکه‌های عصبی انجام شد. در سال ۲۰۰۷ LucasD.Panjer داده کاوی را برای تخمین زمان خطا با مقایسه پنج الگوریتم دسته‌بندی داده‌ها مطرح کرد [۷]. Suma.V و همکارانش به سراغ روش‌های خوشه‌بندی رفته و ثابت کردند که روش K-mean در میان روش‌های مورد بررسی آنها بهینه‌تر بوده است [۸]. CathrinWeiß و همکارانش از روش نزدیکترین همسایه و تشابه جملات با استفاده از موتورارزیابی Lucene^۴ استفاده کردند [۹]. Lucene خود از SVM و یک مدل بولی^۵ برای جستجو و ارزیابی متون استفاده می‌کند. Nagwani نیز در مدل‌های مختلفی از الگوریتم‌های داده‌کاوی متعددی برای اهداف مختلف استفاده کرد [۱۰]، [۱۱]، [۱۲]. اما همه این روش‌ها یا به جنبه معنایی تشابه متون و کلمات توجه نداشتند. یا برای استخراج ضریب تشابه دو متن از الگوریتم‌های مناسب برای این محیط بهره نبردند. الگوریتم‌های استفاده شده بیشتر به تاثیر یک متن و یافتن ضریب تشابه آنها با متن کاربر توجه داشتند، معمولاً استاتیک بودند و هرگز به بهره‌وری روش

^۱ -Bug Tracking

^۲ -GNU product's issue-tracking software

^۳ -Fuzzy Regresion

^۴ -www.lucene.apache.org

^۵ -Boolean model

در محیط با حجم بالا توجه نداشتند. نوآوری کار تحقیقاتی ارائه شده نسبت به مدل‌های قبلی استفاده از یک الگوریتم بهینه با دقت بیشتر در اندازه‌گیری تشابه میان دو متن است به‌طوری‌که با حجم زیاد مستندات و کلمه‌های مشابه در مخازن نرم‌افزاری سازگار باشد. این الگوریتم تشابه جملات و تشابه معنایی بین کلمات را با هم ترکیب می‌کند. به‌طوری‌که، مشابه‌ترین جملات را، با توضیحات داده‌شده برای خطا از سوی کاربر، از میان داده‌های قبلی در مخزن استخراج کند.

۴.۱. فرضیه های تحقیق

۱. مخازن خطای هر نرم افزار منابع ارزشمند دانش هستند.
۲. اطلاعات ثبت شده در مخازن خطا به پیشبرد سریعتر و کم هزینه تر فرآیند تولید نرم افزار کمک می کند
۳. ثبت روش حل مشکل یا خطا در طول فرآیند تولید می تواند به توسعه دهنده پروژه برای حل خطا جدید کمک کند.
۴. برای پیش بینی قسمتی از زمان انجام پروژه و هزینه آن می توان از این اطلاعات استفاده کرد.

۵.۱. اهداف تحقیق

از آنجا که در روش ارائه شده در این تحقیق سعی بر آن شده تا کاستی‌های موجود در روش‌های قبلی به نوعی پوشش داده شود، به دنبال بهترین و مناسب ترین الگوریتم برای جستجو و محاسبه تشابه متون و مستندات در یک مخزن با حجم عظیم داده هستیم. به طوری‌که با بالارفتن این حجم دچار مشکل پیچیدگی محاسبه و بالارفتن زمان پاسخ‌گویی نشود. همچنین تشابه بین دو متن باید شامل تشابه معنایی باشد. به جز تشابه بین کلمات تشابه به صورت مجموعه کلمات. یعنی جملات باید بررسی شود. با توجه به مطالب ذکر شده یافته‌های اصلی این تحقیق را می‌شود در مطالب زیر خلاصه کرد.

۱. ارائه یک روش برای کاوش مخازن خطای نرم‌افزار که کاستی‌های روش‌های قبلی در اندازه‌گیری میزان تشابه بین نمونه و خطا را پوشش دهد.
۲. یافتن خطاهای مشابه با خطای جدید و استفاده از آنها برای تخمین زمان لازم برای رفع خطا
۳. تعیین میزان پیچیدگی خطای جدید با محاسبه میانگین طول عمر خطاهایی با شباهت بیشتر
۴. استفاده از راه کارهای رفع خطا با جستجوی خطاهایی با ضریب تشابه بالاتر با خطای جدید

هدف نهایی کاهش هزینه‌های مالی پروژه‌های نرم‌افزاری و تسريع در فرآیند مهندسی نرم‌افزار است.

۶.۱. نوآوری تحقیق

آنچه بیش از همه به عنوان نوآوری در کار و بهبود روش های گذشته مطرح است:

۱. استفاده از یک الگوریتم معنایی مناسب به همراه تشابه خطی در مرحله اول

۲. استفاده از الگوریتم های خوشه بندی در مرحله دوم.

در این تحقیق برای بالا بردن دقت و سرعت در استفاده از اطلاعات از الگوریتم پیشرفته تر نسبت به کارهای قبلی استفاده شده است. الگوریتم ترکیبی که به جنبه معنایی تشابه نیز توجه داشته باشد. همچنین نتایج نهایی برای بررسی راحت تر و دقیق تر با استفاده از یک روش بهینه خوشه بندی طبقه بندی میشوند، تا دسترسی کاربر به اطلاعات نهایی راحت تر و سریعتر باشد. به بیانی ساده تر اطلاعات نهایی خوانا و جامع باشند.

۷.۱. ساختار پایان نامه

این تحقیق شامل ۵ فصل است که ترتیب فصول و موضوعات مطرح شده به صورت کلی در زیر بیان شده است:

۱.۷.۱. فصل اول: کلیات (همین فصل)

در این فصل یک مقدمه از حوزه این تحقیق ارائه شده و در آن مشخص شده که هدف از این تحقیق چیست و خواهان پاسخ گویی به چه سوالی هستیم. در ادامه همین فصل یک توضیح مختصر از تحقیقات انجام شده در مورد موضوع این تحقیق ارائه شده است. در نهایت نیز یک دید کلی از فصل های مختلف این تحقیق ذکر شده است.

۲.۷.۱. فصل دوم: ادبیات تحقیق

در این فصل به بیان مفاهیم و مقدماتی که این پایان نامه بر اساس آن مطرح شده است، پرداخته می شود. مخازن خطای نرم افزار و انواع آن ها، چرخه عمر یک خطا، همچنین انواع مستندات مرتبط با یک خطا بررسی می شود. پس از آن استراتژی های اندازه گیری تشابه معنایی بین متون و تفاوت آن ها مورد مطالعه قرار می گیرد.

۳,۷,۱. فصل سوم : استفاده از تکنیک‌های داده‌کاوی برای کشف و رفع خطاهای نرم‌افزار بر پایه ضریب تشابه معنایی متن و خوشه‌بندی

ضریب تشابه رشته‌ای یا خطی بر پایه LCS، روش SOC-PMI به‌عنوان یک روش معنایی و تشابه جملات بر اساس عبارات مشترک معرفی می‌شوند. و در نهایت روش اندازه‌گیری شباهت ترکیبی با ادغام و بهینه‌سازی این روش‌ها معرفی می‌شود. روش‌های خوشه‌بندی و محاسبه میانگین داده‌ها مورد نیاز این تحقیق هم بیان می‌شوند. روش پیشنهادی با استفاده از استراتژی‌ها و روش‌های انتخاب‌شده که با محیط داده‌ای مخازن سازگار هستند، برای یافتن خطای مشابه‌تر با خطای جدید تشریح می‌شود و توضیح می‌دهیم که چگونه با یک خوشه‌بندی و محاسبه می‌توانیم میزان پیچیدگی و تخمینی برای زمان رفع خطا ارائه کرد.

۴,۷,۱. فصل چهارم: اعتبارسنجی و توسعه

پیاده‌سازی روش بر روی نمونه داده‌های واقعی که از مخزن خطا استخراج شده با استفاده از برنامه‌نویسی و نرم‌افزار Rapid miner^۱، مقایسه روش از نظر زمان و دقت پاسخ‌گویی با روش‌های دیگر. همچنین روش ارائه شده را با روش‌های قبلی مورد مقایسه و تحلیل موضوعی قرار می‌دهیم.

۵,۷,۱. فصل پنجم : نتیجه‌گیری

در این فصل به جمع‌بندی، نتیجه‌گیری و تبیین دستاوردها و کارهای آینده پرداخته می‌شود. ویژگی‌های این روش با روش‌های دیگر مقایسه شده و زمینه‌های بهبود و نتایج به صورت موردی بیان می‌شود. از آنجا که استفاده از تشابه معنایی در این زمینه در قدم‌های اولیه قرار دارد، در آینده زمینه‌های زیادی برای کار بیشتر وجود دارد که به چند مورد در این فصل اشاره خواهد شد.

^۱-www.apidminer.com

ادبیات تحقیق

۱.۲. فرآیند مهندسی نرم افزار

به زبان ساده فرآیند تولید نرم افزار به مراحل مختلف برای تولید و توسعه یک محصول نرم افزاری انجام می شود اشاره دارد. به صورت کلی فرآیند مهندسی نرم افزار شامل سه مرحله اصلی است: برنامه ریزی (امکان سنجی)، پیاده سازی آزمایش، تست و مستندسازی، استقرار و نگهداری سیستم. آن چه که مد نظر ماست قسمتی از مرحله پیاده سازی، آزمایش، تست و مستند سازی است.

۱.۱.۲. پیاده سازی، آزمایش، تست و مستند سازی

این مرحله در مهندسی نرم افزار بر حسب نوع مدل فرآیند مهندسی شامل قسمت های مختلفی است. در واقع این قسمت تولید نرم افزار است، ساده بگوییم این قسمت برنامه نویسی است. پس از آن نرم افزار باید تست شود. بسیاری از پروژه های بزرگ نرم افزاری مخصوصا پروژه های متن باز در چندین مرحله نسخه های یک نرم افزار را کامل تست می کنند. این کار در شرایط مختلفی انجام می شود که درباره آن توضیح خواهیم داد.

در هر صورت مشکلات نرم افزار باید شناسایی و رفع شوند. مستندسازی نیز در تمام مراحل تولید باید انجام شود. طراحی داخلی نرم افزار برای تعیین اهداف سیستم، نگهداری آینده و ارتقاء و بهبود سیستم هر چند پروژه پایان یافته باشد انجام می شود. همچنین ممکن است این مستندسازی شامل نوشتن ساختار تکه های برنامه، ظاهر برنامه کاربردی داخلی و خارجی هم باشند. این مطلب خیلی مهم است که همه چیز پروژه مستندسازی شود. این مرحله از تولید نرم افزار موضوع تحقیق و راه کار ارائه شده است. بالا بردن بهره وری و پایین آوردن زمان انجام این مرحله از اهداف اصلی این تحقیق هستند.

اگر اهداف مهندسی نرم افزار را موارد زیر در نظر بگیریم. این تحقیق را رسیدن به همه این اهداف را تسهیل می کند.

- افزایش کیفیت، قابلیت اطمینان، قابلیت نگهداری
- رضایت کاربران و سهامداران
- کاهش هزینه های جانبی و پشتیبانی
- تحویل به موقع
- استفاده از مولفه های استاندارد
- استفاده مجدد

مرحله مستندسازی یا به روایتی ثبت تمامی اطلاعات برآمده از پروژه، ارتباط تنگاتنگی با مرحله تست یا رفع خطا دارد. همچنین هر دو این مراحل وابسته به این تحقیق هستند. اینکه چگونه می توان از اطلاعات ذخیره شده به شکل های مختلف در طول عمر پروژه برای پیشبرد مرحله تست استفاده کرد، نیازمند راه کارهایی هوشمند در حوزه داده کاوی است. قبل از هر چیز باید بدانیم که در یک پروژه نرم افزاری چگونه مستندات و اطلاعات متنی ذخیره می شود. منظور از خطا در این تحقیق چیست؟ در ادامه انواع مخازن داده و اطلاعات یک پروژه نرم افزاری معرفی می شود.

۲.۲. انواع مخازن داده

۱.۲.۲. کد اصلی:

کد اصلی بخش قابل اجرا و رفتار یک توسعه نرم افزاری است. که در نهایت به صورت فرمت اجرایی به مشتری تحویل داده می شود. که عموماً به عنوان مهمترین داده از سوی توسعه دهندگان مورد توجه قرار می گیرد. مخزن حاوی این اطلاعات شامل تعدادی از منابع کد و اسناد در یک یا چند زبان مختلف برنامه نویسی است. این اسناد معمولاً به موجودیت های منطقی به نام های ماژول^۱ یا بسته گروه بندی می شوند. تمام این مجموعه اطلاعات کد اصلی سیستم نامیده می شود. برای کاوش این متون تمرکز روی شناسه ها (متغیر، نام)، توضیحات و رشته های اصلی داخل کد اصلی است. معمولاً کلمات کلیدی و نمادها در نظر گرفته نمی شوند.

۲.۲.۲. مخازن خطا (سیستم ردیابی خطا BTS)

این مخازن برای ذخیره اطلاعات مربوط به ایجاد و حل خطا، مشخصات ارتقاء سیستم و کلیه اقدامات دیگر در مرحله تعمیر و نگهداری استفاده می شوند. معمولاً هنگامی که توسعه دهندگان کاربران به مشکل یا خطایی در یک سیستم نرم افزاری مواجه می شوند، یادداشتی درباره این خطا در پایگاه داده خطا در موضوع مربوطه ذخیره می شود. این اطلاعات شامل: علت و مکان وقوع خطا در برنامه و اینکه چگونه خطا باعث ایجاد اشکال و خلل در روند برنامه شده است. پس از آن یک یا چند متخصص، موضوع ایجاد شده را برای رفع مشکل بررسی می کنند. چنانچه خطا برطرف شود موضوع در فرم مربوطه بسته می شود. تمام این اطلاعات در مخازن و پایگاه های خطا ذخیره می شوند. عمومی ترین سیستم های مخازن خطا Bugzilla، Trac هستند.

^۱-Module

^۲-bug-tracking system

اگرچه تا به امروز سیستم‌های متعددی ساخته شده‌اند. در حالت عادی بین خطا^۱، نقص^۲، عیب^۳ تفاوت قائل می‌شویم، اما در این تحقیق همه را با لفظ خطا و هم معنی در نظر می‌گیریم.

۳،۲،۲. لیست نامه‌ها و گفتگوهای ثبت شده

لیست ایمیل‌ها (یا آرشیو بحث‌ها) همراه با گفتگوهای ثبت شده بین افراد دخیل در یک پروژه آرشیوی از ارتباطات متنی توسعه‌دهندگان، مدیران و ذینفعان آن پروژه هستند. لیست متنی متشکل از بسته‌های الکترونیکی که شامل سه قسمت:

۱. سرآیند (فرستنده، گیرنده و زمان ارسال)

۲. بدنه پیغام (متن داخل ایمیل)

۳. مجموعه‌ای از فایل‌های پیوست‌شده (مستندات اضافی که همراه ایمیل فرستاده می‌شود) می‌باشد.

شرح گفتگوها شامل ثبت مکالمات فوری بین ذینفعان پروژه، که بر حسب زمان یا نویسنده دسته‌بندی شده‌اند، می‌باشد.

۴،۲،۲. پایگاه داده کنترل منبع (پایگاه داده کنترل ویرایش‌ها)

سیستمی برای ثبت تاریخ تغییرات (ویرایش‌ها) به همراه خود ویرایش و اطلاعات دیگر به صورت اسناد و اطلاعات متنی است. توسعه‌دهندگان معمولاً تاریخ و زمان ویرایش یک کد اصلی را در پایگاه داده‌هایی ذخیره می‌کنند. پایگاه داده‌های کنترل کد رایج مانند cvs [۱۳] و svn [۱۴]، به توسعه‌دهندگان اجازه می‌دهند به یک کپی از مخزن سراسری و جهانی، در سیستم فایل‌های محلی خود، دسترسی داشته باشند. اسناد موجود را ویرایش کنند، یا اطلاعاتی اضافه یا کم کنند و یا ساختار دایرکتوری این مخازن را تغییر دهند. همچنین می‌تواند در مخزن اصلی سند یا اطلاعات جدید محلی ایجاد کند.

بنابراین کنترل بازبینی‌ها دو نتیجه مهم در بر خواهد داشت:

- اول اینکه به توسعه‌دهندگان اجازه می‌دهد، مستقل از کسانی که به مخازن دسترسی دارند، فایل‌های روی سیستم‌های خود را تغییر دهند. پس از آن که تغییرات تایید شده ایجاد شد بقیه می‌توانند این تغییرات را بررسی کنند. این استقلال کاری اجازه می‌دهد که یک چرخه کار موازی بدون نیاز به ارسال ایمیل و گفتگو و نیز بدون تغییرات ورژن برنامه به عقب و جلو تشکیل شود.

^۱ -bug

^۲ -defect

^۳ -fault

- دوم اینکه زمان و تاریخ همه اطلاعات و مستندات به صورت خودکار ثبت و نگهداری می‌شود. اگر نسخه‌های قبل نرم‌افزار نیاز بود، توسعه‌دهندگان به راحتی می‌توانند به نسخه‌های قبل سیستم دسترسی داشته باشند و سیستم را به نسخه قبلی برگردانند.

۵.۲.۲. اطلاعات طراحی و نیازمندی‌های سیستم

مستندات نیازمندی‌ها، معمولاً در ارتباط با مشتری و یا با تاییدهای او تنظیم می‌شود. این اسناد لیستی از نیازهای مشتری است که خواهان انجام آن توسط سیستم است. این نیازها به دو صورت دسته‌بندی می‌شوند. اینکه چه نیازهایی را سیستم باید برطرف کند و چگونه و با چه کیفیتی موردانتظار مشتری است. اطلاعات طراحی نیز شامل تمام اطلاعات مربوط به طراحی معماری و الگوریتم‌های مهم و مورد استفاده^۱ سیستم است. طراحی سیستم می‌تواند به شکل نمودار (مانند UML) و یا به صورت متون جریان کار نمایش داده شوند.

۶.۲.۲. مخازن شرح اجرا

اطلاعات مربوط به خروجی یک برنامه در حین یک یا چند بار اجرا، بر حسب آن چه که در قسمت تست مشخص شده ثبت می‌شود. این اطلاعات شامل لیستی از زمان و دلیل اجرا شدن متدهای یک سیستم، مقادیر قطعی متغیرها و جزئیات مربوط به شرایط اجراست. این لیست هنگامی مفید است که فرآیند اشکال‌زدایی در مقیاس وسیع همزمان با هزاران یا میلیون‌ها فرآیند دیگر از سوی افراد مختلف انجام شود. در این هنگام پیدا کردن اشکالات فردی لازم اما کاری دشوار و زمان بر است. این مخازن کمک می‌کنند به تک تک این اجراها با جزئیات دسترسی داشته باشیم.

۷.۲.۲. مخازن سیستم نرم افزار

مخازن سیستم نرم‌افزار شامل مجموعه‌ای از سیستم‌های مختلف به همراه کد اصلی آنها که افراد علاقه‌مند می‌توانند به راحتی در آن جستجو کرده و از آنها استفاده نمایند. از مخازن مهم رایج در این حوزه می‌توان به Source Forge [۱۵] و Google Code [۱۶] اشاره کرد. این مخازن شامل آرایه وسیعی از اطلاعات در حوزه‌های مختلف پروژه‌های نرم‌افزاری است که می‌توان اطلاعات غنی را در فازهای مختلف آن استخراج کرد. آنچه به عنوان منابع داده در این تحقیق مورد استفاده قرار می‌گیرد همان مخازن خطا است. که با استفاده از یک سیستم ردیابی خطا نمونه‌های آماری لازم استخراج می‌شود [۱۷].

^۱ -Use Case

۳.۲. خطای نرم‌افزاری

اشکال نرم‌افزاری، خطا، نقص و یا شکست ایجاد شده در برنامه کامپیوتری است، که باعث تولید یک نتیجه نادرست یا غیرمنتظره می‌شود. یا باعث بروز رفتار ناخواسته از سیستم می‌گردد. رفع این مشکلات در پروژه‌های نرم‌افزاری از سخت‌ترین و پرهزینه‌ترین مراحل مهندسی نرم‌افزار است. این خطاها می‌تواند در پروژه‌های بزرگ و مخصوصاً پروژه‌های متن‌باز^۱ از سوی همه کاربران و توسعه‌دهندگان گزارش شوند. از این رو نرم‌افزارهای ردیابی خطا به کمک آمده تا بتوان این گزارشات را ثبت و پی‌گیری کرد. خاطر نشان می‌کنم که در این تحقیق همان‌طور که قبلاً گفتیم مشکلات اعم از نقص، شکست یا خطا با عنوان خطای نرم‌افزاری معرفی می‌شود.

۱.۳.۲. سیستم ردیابی خطا

سیستم ردیابی خطا یک برنامه نرم‌افزاری است که برای کمک در سهولت پی‌گیری خطاهای نرم‌افزاری در طول عمر پروژه طراحی و ارائه شده است. این سیستم یک نرم‌افزار کاربردی برای تیم توسعه است به‌طوری که به آن‌ها در تضمین کیفیت و کاهش هزینه‌ها و پیشرفت سریع پروژه کمک می‌کند. این سیستم‌ها به کلیه افراد دخیل در یک پروژه اعم از کاربران، توسعه‌دهندگان و ذینفعان دیگر اجازه می‌دهد که خطای رویت شده در نرم‌افزار را در هر مرحله از روند تولید که باشد گزارش داده و پیگیری کنند. بسیاری از پروژه‌های بزرگ از این نرم‌افزارها استفاده کرده و پروژه‌های خود را در نسخه‌های مختلف ارائه می‌دهند. تا در مراحل مختلف از ارائه نسخه جدید خطاهای نسخه قبلی که دیگران گزارش داده‌اند پیگیری و رفع کنند. پروژه در هر حالتی ارائه شود، چه به صورت مرحله‌ای یا کامل، متن‌باز یا بسته و با هر متد استفاده شده در مهندسی نرم‌افزار، یک سیستم ردیابی خطا در پیشبرد بهتر و سریع‌تر فرآیند مهندسی نرم‌افزار بسیار ارزشمند است.

بسیاری از این نرم‌افزارها برای استفاده عموم طراحی شده و بقیه تنها برای استفاده در یک پروژه طراحی می‌شوند. معمولاً سیستم‌های ردیابی خطا با دیگر اپلیکیشن‌های مدیریت پروژه یکپارچه و مجتمع می‌شوند [۱۸]. مهمترین قسمت این سیستم‌ها مخازن اطلاعاتی آن‌هاست که کلیه اطلاعات مربوط به خطا را در خود ذخیره می‌کنند. باید دید در یک سیستم ردیابی خطا ایده‌آل چه اتفاقی می‌افتد و یک خطا چه چرخه‌ای را در طول عمر خود در آن طی می‌کند. به عنوان مثال باگزیلا^۲ یک سیستم ردیابی خطا است که در سال ۱۹۹۸ توسط تری ویسمن^۳ برای شرکت موزیلا^۴ نوشته شد.

^۱ -open source

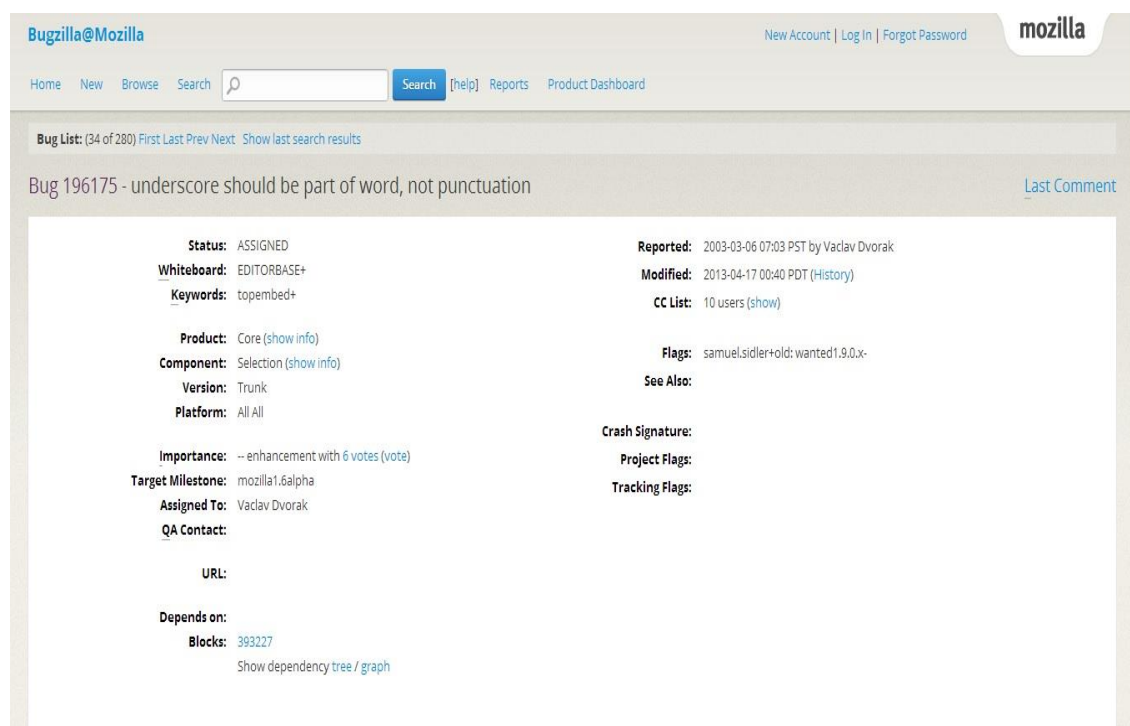
^۲ -Bugzilla

^۳ -Terry Weissman

^۴ -www.mozilla.org

در ابتدا باگزیلا به عنوان یک سیستم ردیابی خطا متن باز از داخل خانه برای کاربران مرتبط با نت اسکپ^۱ طراحی شد. ابتدا به زبان TCL^۲ و سپس به Perl برای راحتی و سهولت استفاده و دخالت همگان در پیشرفت کار برگردانده شد. این سیستم برای استفاده در پروژه‌های متن باز بهینه‌سازی شد. به طوری که شرکت‌ها و پروژه‌های بزرگی چون یاهو، ناسا، Gnome، KDE، Apache، RedHat، و یکی مدیا و خود موزیلا از آن استفاده کردند هر خطا یک شناسه (ID)، عنوان، تاریخ، وضعیت، نام محصول و مولفه، میزان شدت سختی و پیچیدگی، اولویت و غیره را در خصوصیات خود دارد [۹] (شکل ۲-۱).

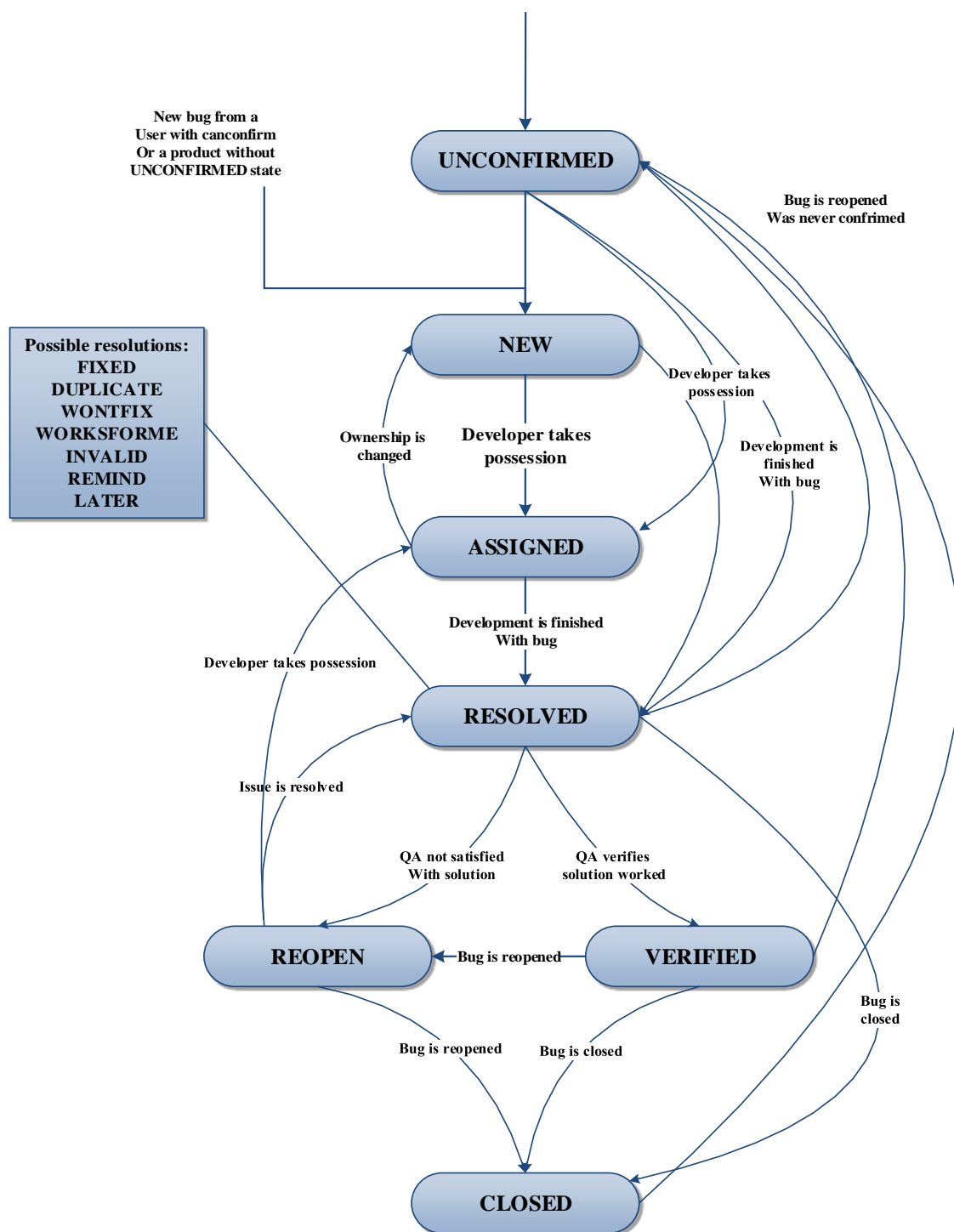
شکل ۲-۱- مشخصات یک موضوع (خطا) در نرم افزار Bugzilla



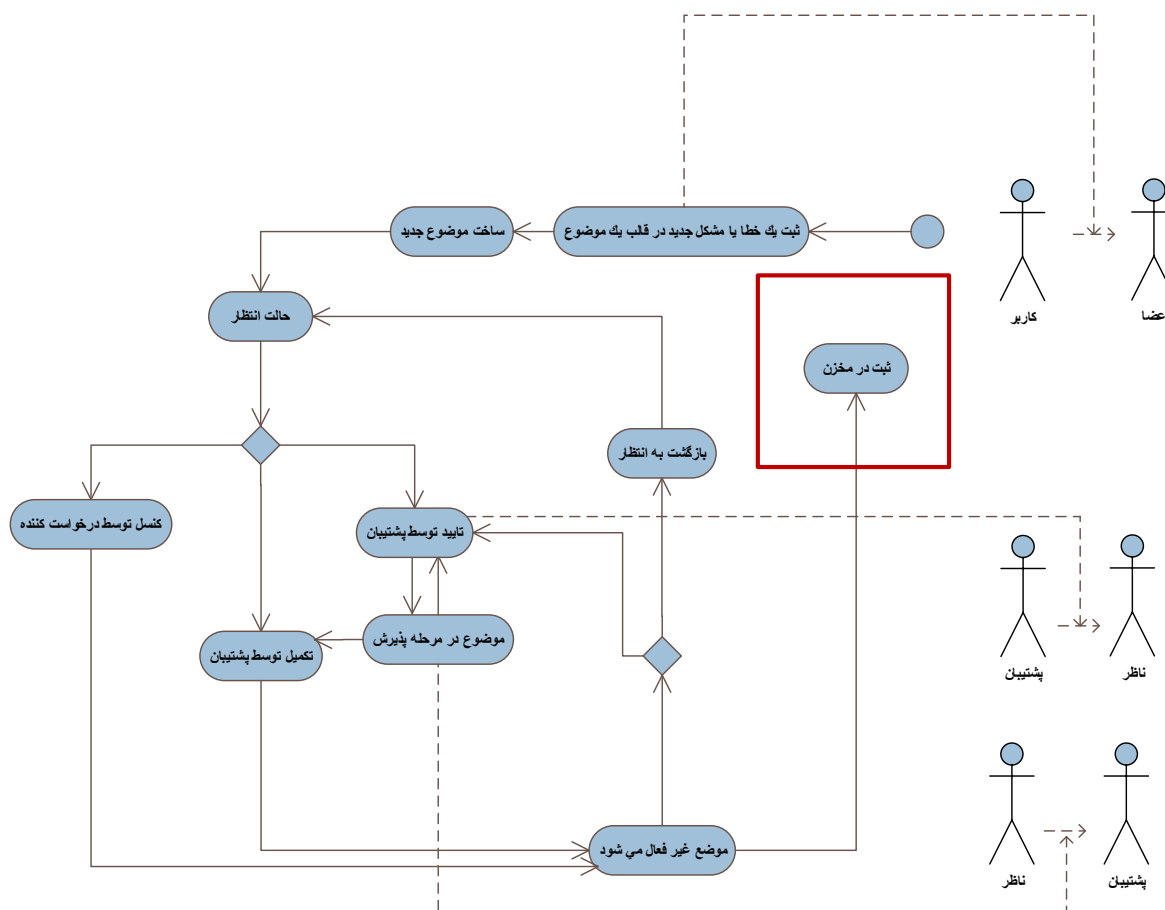
یک خطا از زمان ایجاد تا زمان بسته شدن (حل مشکل) چرخه‌ای را در سیستم دارد (شکل ۲-۲) نمودار فعالیت برای یک مخزن خطا در شکل ۲-۳ نشان داده شده است.

^۱-Net Scape

^۲-Tool Command Language



شکل ۲-۲ - چرخه یک خطا در یک سیستم مدیریت خطا [۱۹]



شكل ۲-۳ - نمودار جریان کار^۱ یک سیستم ردیابی خطا

۱- خطا جدید (NEW) یا تایید شده وارد سیستم می شود (unconfirmed) اگر اعتبار خطا توسط یکی از اعضای تیم یا هسته مدیریتی تایید شد، به عنوان یک گروه جدید ثبت می شود.

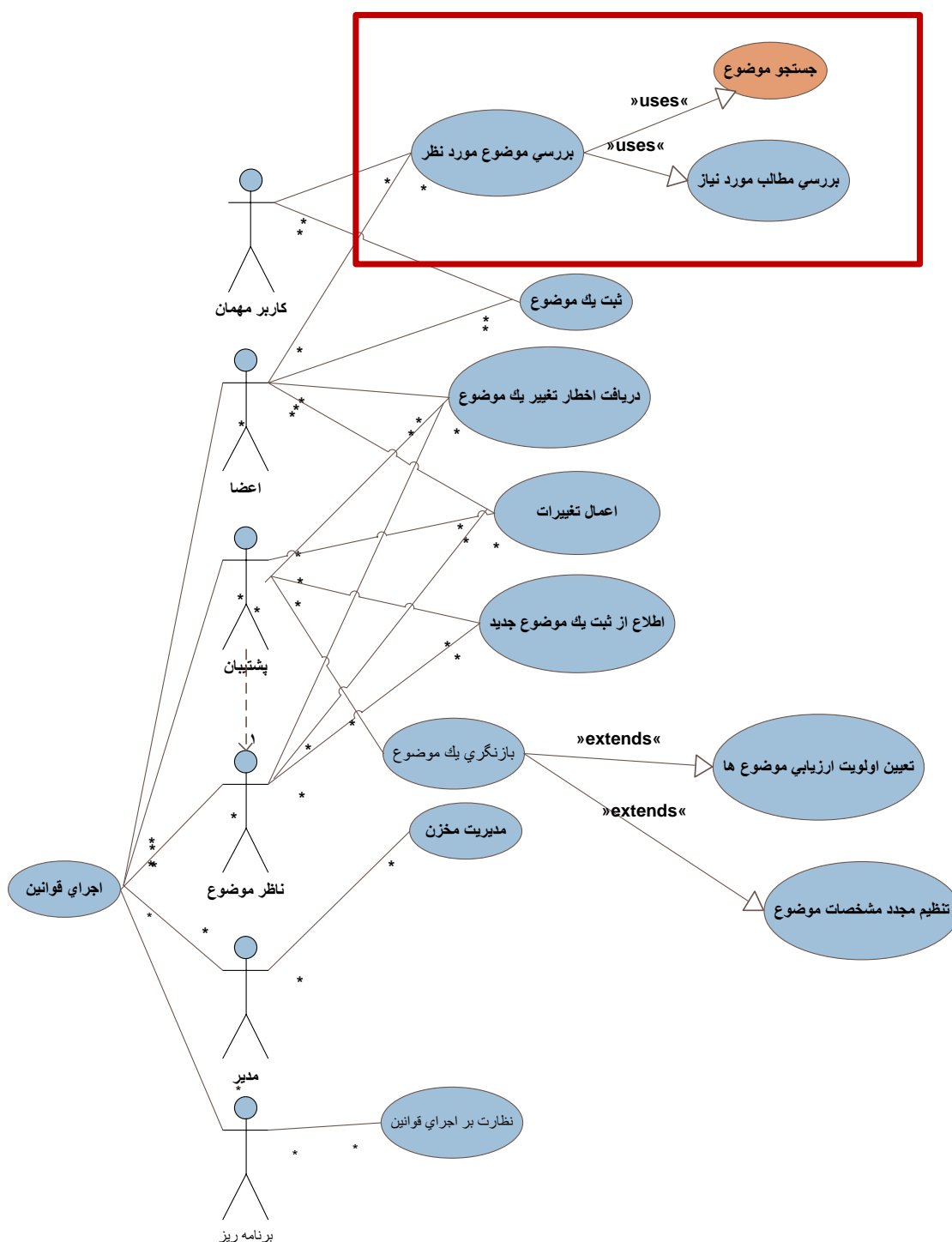
۲- مدیران و کسانی که مسئول ارزیابی خطاها هستند باید وجود خطا را تایید کنند. اعتبار و غیرتکراری بودن آن را نیز بررسی کنند. اگر خطا تایید شد به قسمت تایید شده ها (Assigned) رفته، تا در اختیار تیم توسعه برای حل و فصل و رفع خطا قرار گیرد. در غیر این صورت به قسمت تایید نشده (unconfirmed) برای بررسی بیشتر می رود.

۳- اگر خطا رفع شد برچسب حل شده (Resolved) خورده و آن را به وضعیت حل شده انتقال می دهند.

^۱ -Work Flow Diagram

۴- اعضای تیم ممکن است یک خطا حل شده را به یکی از حالت‌های سازگار (Verified) و سپس بسته (Closed) اتصال دهند. یا دوباره به شکلی در رابطه با این خطا مواجه شدند و آن را به حالت دوباره باز یا تایید شده بفرستند.

خطاهایی که به مرحله Resolved بازگردانده می‌شوند یا نامعتبرند یا تکراری هستند و نیاز به بررسی مجدد دارند و یا به‌طور کامل رفع شده‌اند. چرخه عمر یک خطا از مرحله‌ی جدید تا مرحله خروج از حل تعریف می‌شود. شکل ۲-۴ وظایف هر یک از اشخاص در ارتباط با مخزن خطا را نشان می‌دهد. دو کادر قرمز در دو نمودار شکل ۲-۴ و شکل ۲-۳ محدوده کاری این تحقیق را نشان می‌دهند.



شکل ۲-۴- نمودار مورد کاربرد^۱ یک سیستم ردیابی خطا

همان طور که در شکل ۲-۵ می بینید، یک خطا دارای مشخصات متنی است. برچسب، وضعیت خلاصه و در قسمت جزئیات توضیحات و راه حل های ارائه شده با ذکر نام نویسنده و تاریخ برای هر خطا درج شده است. همین اطلاعات و متون ثبت شده دامنه کار تحقیقاتی ماست. این داده ها منبع غنی اطلاعات و دانش سودمند است. اینکه چگونه این داده ها را کاوش کنیم و چه اطلاعاتی در آنها

¹ - Use Case Diagram

نهفته است، گام‌های تحقیق ما را می‌سازند. از آنجایی که این داده‌ها توسط افراد مختلف و با استفاده از قوانین نحوی صحیح و ناصحیح و گاه کلمات مشابه در یک حوزه استفاده شده‌اند نیاز به یک جستجو معنایی به شدت احساس می‌شوند.

ID	Status	Version	Component	Severity	Summary
1026187	NEW	unspecified	General	normal	Warnings displayed when running tests
1026149	NEW	unspecified	General	normal	combining multiple ui buttons
1025939	NEW	unspecified	General	normal	TEST-UNEXPECTED-FAIL tests/test-native-options.testSimplePrefs [Exception... "Component (NS_ERROR_UNEXPECTED) [nsIPrefBranch.setBoolPref]
1025356	UNCONFIRMED	unspecified	General	normal	Firefox 30 extension with Add-on SDK throwing NS_ERROR_FAILURE in sdk/timers.js
1025315	NEW	unspecified	General	normal	only run sdk/preferences/native-options.enable() if optionsURL is data:text/xml,<placeholder/>

شکل ۲-۵- مشخصات یک خطا

پایه و اساس این تحقیق جستجو جملات مشابه معنایی است. به‌صورتی که آن‌دسته از جملات که دارای تشابه معنایی بیشتری با جمله یا کلمه کاندیدای ما هستند، جداسازی شوند. چگونگی این کار در سه مرحله خلاصه می‌شود: ابتدا با استفاده از تکنیک‌هایی که در ادامه پیشنهاد می‌شود از میان داده‌های آماری که از یک مخزن خطا انتخاب شده، میزان تشابه هر خطا (بر اساس داده‌های متنی مانند توضیحات، سرآمد، راه‌حل و غیره) با خطای جدید اندازه‌گیری می‌شود. در مرحله بعد این خطا با استفاده از یک الگوریتم خوشه‌بندی که سعی شده یک الگوریتم بهینه باشد به دسته‌های جدا بر اساس ضریب تشابه محاسبه شده تقسیم می‌شوند.

این کار به کاربر کمک می‌کند که نمونه‌های مشابه گزینش شده را در خوشه‌های طبقه‌بندی شده از نظر سختی و پیچیدگی کار در اختیار داشته‌باشد. گاهی ممکن است کاربر خود بر اساس تجربه حدس بزند که مشکل زیاد پیچیده نیست. یا برعکس مشکل به زمان زیادی برای حل نیاز داشته باشد، در این حالت این خوشه‌بندی می‌تواند در کاستن زمان یافتن راه‌حل کمک قابل توجهی داشته باشد. آنچه این کار را از کارهای انجام شده قبلی متمایز می‌کند توجه به بهره‌وری و بهبود کارایی است. یک مخزن خطا مجموعه‌ای گسترده و بزرگ از داده‌هاست که هر روز بزرگتر می‌شود. با بزرگ شدن نمونه آماری و نیاز به دقت در کار ابعاد پیچیدگی این عمل روشن خواهد شد. هدف اصلی این کار استفاده از تکنیکی است که دقت بالاتری داشته باشد. جستجو در یک بانک داده که متون در آن از قانون خاص برای رعایت اصول تحریر و قواعد نحوی استفاده نشده است، و کسانی که داده‌ها را ثبت کرده‌اند ممکن است از کلمات مشابه زیادی استفاده کرده باشند، نیاز بیشتری به یک جستجوی معنایی دارد.

۴.۲. تحقیقات پیشین در حوزه داده‌کاوی در مخازن خطا

در سال ۱۹۹۲ اولین نرم‌افزار ردیابی خطا GNATS^۱ شروعی برای MSR بود. با ارائه اولین نرم‌افزار ردیابی خطا، مدیریت و استفاده از دانش نهفته در خطاها اولین قدم‌های خود را برداشت. در این مدت محققان سعی کردند برای استخراج بهینه و مفیدتر دانش، مدل‌ها و روش‌های جدیدی با استفاده از الگوریتم‌های مختلف ارائه کنند. بیشتر این روش‌ها از الگوریتم‌های دسته‌بندی برای دسته‌بندی اطلاعات و استفاده از دانش آن‌ها به صورت طبقه‌بندی شده استفاده کردند. ابتدا مروری کوتاه بر تحقیقات قبلی و روش‌های بررسی شده خواهیم داشت. نواقص و کمبودهایی که به موجب آن‌ها این تحقیق انجام شده و روش ارائه شده‌است نیز بررسی خواهد شد.

سال ۲۰۰۰، Junzo Watada به منظور برآورد تعداد خطاها در یک پروژه نرم‌افزاری از رگرسیون فازی استفاده کرد [۶]. وی مجموعه سئوالاتی را برای استفاده در سیستم فازی در مراحل مختلف از یک پروژه مطرح کرد. در این روش تمام داده‌ها برای جستجو هدف قرار نمی‌گیرد و این خود ضعف بزرگی است.

Lucas D. Panger در سال ۲۰۰۷ مقایسه‌ای برای دسته‌بندی مخازن خطا با استفاده از پنج الگوریتم دسته‌بندی 1-R، 0-R، درخت تصمیم‌گیری C4.5، Naïve Bayes، و رگرسیون لجستیک، در داده‌های گرفته شده از Bugzilla انجام داده‌است [۵]. این کار با محاسبه درصد خطاهایی که دسته‌بندی شده‌اند و ضریب Kappa برای هر کدام با استفاده از نرم‌افزار Weka انجام داده است. داده‌ها بر اساس طول عمر و توضیحات متنی دسته‌بندی شده‌اند. این تحقیق تنها مقایسه‌ای بین الگوریتم‌های مختلف در یک دسته‌بندی ساده برای داده‌های موجود در یک مخزن است، که کمتر به شباهت بین یک خطا جدید و خطاهای دیگر پرداخته شده است.

در همان سال Cathrin Weib و همکارانش با استفاده از داده‌های گرفته شده از JBoss در سه مرحله طول عمر خطای جدید را تخمین زدند [۹]. ابتدا به کمک الگوریتم نزدیک‌ترین همسایه α -KNN گزارش‌های مشابه را از منبع استخراج می‌کند و در مرحله بعد با استفاده از موتور جستجو Lucene داده‌ها با شباهت متنی بیشتر استخراج می‌شود. Lucene از SVM و یک مدل بولی^۲ برای جستجو و ارزیابی متون استفاده می‌کند. در این تحقیق اگر موضوع با شباهت با توضیحات خطای جدید پیدا نشد، باز الگوریتم نزدیک‌ترین همسایه α -KNN را برای استخراج داده استفاده می‌شود پس از این مراحل داده‌های فیلتر شده می‌تواند در تخمین طول عمر خطا جدید و حل آن کمک می‌کند. این روش جزء روش‌های ابتدایی با به کارگیری اندازه‌گیری تشابه بین متون بود. اما موتور و الگوریتم آن نه تنها به معنای کلمات توجه نمی‌کرد بلکه جزء روش‌های ابتدایی جستجو در متون متشابه بود.

^۱-GNU product's issue-tracking software

^۲-Boolean model

Nagwani در سال ۲۰۰۹ روشی را با دامنه خطاهای حوزه GUI ارائه کرد. در این تحقیق خطاهای مربوط به GUI را از پروژه مختلف را به عنوان داده در نظر گرفت. از نظر او خطا در طول فرآیندی رخ می دهد که این فرآیند شامل چند رویداد متوالی است [۱]. مدل وی در سه مرحله داده ها را فیلتر می کند. در مرحله اول داده ها بر اساس رویداد، در مرحله بعد برای قسمت های مختلف هر رویداد مقایسه می شوند. در مرحله آخر براساس ملزومات هر قسمت تفکیک و با هم مقایسه می شوند. در این روش خطاهایی که از نظر مرحله و رویدادهای موجود در فرآیند دارای تشابه بیشتر با خطای جدید هستند از منابع استخراج می شود. این روش اگرچه به مفاهیم و معانی توضیحات هر موضوع خطا بیشتر توجه کرده، اما در پیاده سازی موانع زیادی وجود دارد. مراحل توسط نیروی انسانی مجزا و مشخص شوند و تمام این رویدادها در طول یک فرآیند خطا دار باید مشخص باشند. همچنین محدود به خطاهایی از که در یک فرآیند مشخص در GUI است.

یک سال بعد Nagwani و Bhansali با استفاده از الگوریتم K-means مجموعه ای از داده ها شامل طول عمر و شماره داده را به خوشه هایی تقسیم کردند، تا بتوانند میزان پیچیدگی یک خطا را با تخمین ذهنی طول عمر یک خطای جدید را مشخص کند. به این ترتیب ابتدا کاربر برای طول عمر خطا جدید تخمینی می زند و این تخمین با مراکز خوشه های ایجاد شده بر اساس طول عمر داده ها، پیچیدگی خطا جدید را مشخص می کند [۱۱]. این روش دو مورد مهم توجه نشده اول اینکه تخمین ذهنی طول عمر یک خطا زیاد اصولی و دقیق نیست. دوم اینکه تعیین میزان پیچیدگی یک خطا با یک پس زمینه ذهنی کمک چندانی در پیشرفت پروژه و بهره وری در هزینه و زمان آن ندارد.

در تحقیقی دیگر در همان سال، وی با استفاده از سه الگوریتم تشابه جاکارد، کسینوسی و TF-IDF میانگینی از تشابه موجود بین قسمت های متنی مختلف از خطای جدید و خطاهای دیگر در منبع خطا را محاسبه می کند [۱۰]. نتایج را برای هر الگوریتم با هم مقایسه کرده است. این کار با محاسبه تشابه بین بخش موضوع، توضیحات و دیگر بخش های متنی از هر موضوع موجود در منبع خطا را با خطای جدید، به صورت جداگانه محاسبه می کند. این نتایج با ضرایب وزنی مشابه با هم جمع می شوند و به عنوان میزان شباهت خطای جدید با خطا داخل مخزن در نظر گرفته می شود. برنامه ای که Nagwani ارائه کرده است این کار را با استفاده از هر سه الگوریتم انجام می دهد.

برای نمونه در این تحقیق، با استفاده از ضریب کسینوسی برای یک نمونه جدید، نتایج محاسبه و چهار خطا با شباهت بیشتر، به عنوان کمک برای حل خطا و محاسبه میانگین طول عمر آنها، به عنوان تخمینی برای طول عمر خطا استخراج شده است. در این روش برای بررسی میزان تشابه بین متون خطا جدید و داخل مخزن از روش های خطی استفاده شده که برای شروع مناسب است، اما کافی نیست. این الگوریتم ها معایبی دارند که الگوریتم های بهتری را می توان استفاده کرد.

به عنوان مثال الگوریتم‌های معنایی که در آنها علاوه بر تشابه لغوی و نحوی تشابه معنایی بین جملات و متون را در نظر گرفته می‌شود.

در سال ۲۰۱۱، Nagwani تحقیق خود را روی نویسندگان و توسعه‌دهندگان ثبت شده در مخزن خطا و موضوعاتی که جوابگوی آنها بوده‌اند یا نظرات موثری در رفع خطا داشته‌اند، متمرکز کردند [۱۲]. به این صورت که روی تعدادی داده که از Bugzilla استخراج شده بود کلمات تکراری مهم در این موضوع‌ها را به عنوان داده استخراج می‌کند. لیستی از توسعه‌دهندگان و کاربران ثبت کننده به تفکیک کلمات اصلی تکرار شده لیست می‌شود. در این حالت می‌توان کسانی را که در حوزه موضوع خطای جدید کار کرده‌اند، برای کمک به حل مشکل شناسایی کرد. در این روش تشخیص موضوع اصلی و کلمات مهم باید مطابق با مجموعه و در مقایسه با آن استخراج شود. انتخاب و یا تشکیل مجموعه بهینه خود نیازمند ارائه روشی مناسب است.

Suma.V و همکارانش در سال ۲۰۱۲ به استفاده از چند روش خوشه‌بندی و مقایسه آن‌ها اکتفا کرد. آنها هشت الگوریتم را برای دسته‌بندی داده‌ها در مخازن خطا استفاده کردند که در میان آنها K-means در برابر افزایش داده‌ها نتایج بهینه‌تری در اختیار کاربر می‌گذارد [۸]. الگوریتم‌ها با استفاده از نرم‌افزار Weka پیاده‌سازی شده‌اند. این روش نیز تنها به دسته‌بندی ساده اکتفا شده و به نوع داده‌ها که متنی هستند توجه نشده است.

Nagwani در جدیدترین تحقیق خود در ۲۰۱۳ الگوریتم‌های NB، J84، SVM، CC و CLUBAS را روی داده‌های گرفته شده از پروژه Mozilla برای دسته‌بندی استفاده کرد و نتایج را تحلیل کرده است [۲۰]. به عنوان نتیجه نشان داده است که NB و J84 با زیاد شدن حجم داده‌ها افت دقت ندارد.

پس از مطالعه‌ای اجمالی روی BTS ها و روش‌هایی که تاکنون برای جستجو و استفاده دانش موجود در این سیستم‌ها نشان می‌دهد، که سیر پیشرفت آن مدتی است که متوقف شده. در واقع پیشرفتی چشم‌گیر در آن‌ها دیده نمی‌شود. آنچه نبودش بیش از همه در کاوش در داده‌های مهم موجود در این مخازن، به چشم می‌آید، هوشمند سازی این کاوش هاست. داده‌های موجود در این مخازن مجموعه‌ای از متون و ارقام است. پس جستجو و کاوش پیشرفته با در نظر گرفتن جنبه‌های معنایی موجود در آنها، مطمئناً نتایج بهتری را در اختیار کاربران قرار خواهد داد.

در این تحقیق سعی شده به این جنبه از داده کاوی مخازن، یعنی استفاده از الگوریتم‌هایی که به مفهوم و معنای داده‌ها و استخراج دانشی بهینه و مختصرتر توجه شود.

۵.۲. اندازه گیری شباهت بین متون

پیدا کردن شباهت بین دو متن و جمله با پیدا کردن شباهت بین کلمات آغاز می شود. کلمات به صورت معنایی یا لغوی می توانند به هم شبیه باشند. شباهت لغوی به این معنی است که حروف کلمات با هم مشابه باشند. و شباهت معنایی حالتی است که دو کلمه در یک زمینه و به یک نوع و معنی استفاده شوند.

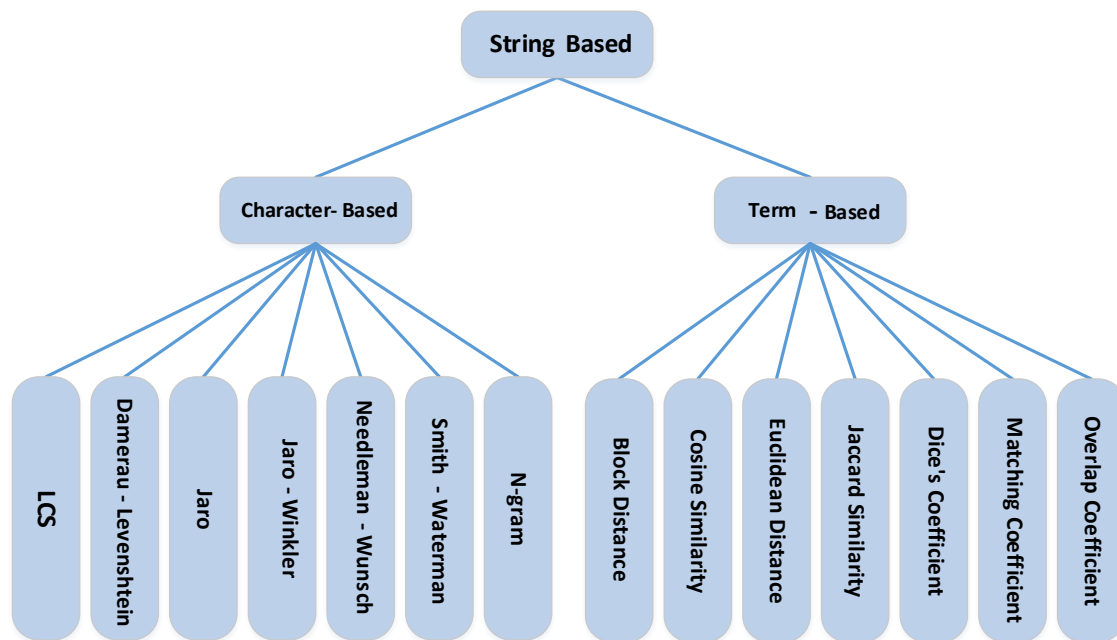
شباهت معنایی در روش ها مقوله ای جدا از شباهت خطی است. شباهت معنایی به صورت شباهت بر پایه مجموعه و شباهت بر پایه دانش معرفی می شود. شباهت خطی کاملاً متریک بوده و تشابه دو جمله را بر اساس توالی دستهای و ترکیب کاراکترها و کلمات داخل رشته بررسی می کند. شباهت معنایی، شباهت بین کلمات دو جمله براساس اطلاعات مجموعه ای که شامل آن می باشد، را مشخص می کند. ابتدا متداول ترین روش های هر دو معیار را به اختصار معرفی کرده و سپس روش مورد استفاده این تحقیق را با شرح کامل دلایل انتخاب شرح می دهیم.

۵.۲.۱. شباهت خطی

معیار شباهت خطی بر پایه توالی و موقعیت کاراکترها بنا شده است. در اندازه گیری این شباهت میزان (فاصله) شباهت دو جمله بر اساس همین دو معیار یعنی توالی خطی و موقعیت کلمه در جمله اندازه گیری می شود. از چهارده روش و معیار اندازه گیری ارائه شده در طول تحقیقات عمومی، در این گروه هفت معیار مربوط به گروه توالی خطی^۱ و هفت معیار و روش به گروه موقعیت کلمه^۲ تعلق دارند. (شکل ۲-۶)

^۱-Term-based

^۲-Character-based



شکل ۲-۶ - انواع الگوریتم های تشابه خطی

۱،۱،۵،۲. اندازه گیری شباهت بر پایه کاراکتر^۱

▪ بزرگترین زیر دنباله مشترک^۲ (LCS)

روش‌های مختلفی برای پیدا کردن بزرگترین زیر دنباله در مجموعه‌ای از دنباله‌ها. معیار شباهت بزرگترین رشته پیوسته از کاراکترهاست که در هر دو جمله موجود باشند. هدف در این روش مقایسه دو رشته S_1 و S_2 است به طوری که رشته S_3 با حفظ ترتیب، و نه لزوماً متوالی در هر رشته موجود باشد.

▪ Damerau-Levenshtein

تعریف فاصله بین دو رشته با شمارش حداقل تعداد عملیات مورد نیاز برای تبدیل یک رشته به دیگری که در آن عملیات باید درج، حذف و یا جایگزینی یک واحد یا جابجایی در حرف مجاور باشد.

▪ Jaro و Jaro-Winkler

بر پایه تعداد و ترتیب کاراکترهای مشترک بین دو رشته تعریف می‌شود به طوری که فاصله دو کلمه مشترک با توالی یکسان در دو خط هر چه کوتاه‌تر باشد دو جمله شباهت بیشتری دارند.

$$d_j = \begin{cases} 0 & \text{if } m = 0 \\ \frac{1}{3} \left(\frac{m}{|s_1|} + \frac{m}{|s_2|} + \frac{m-t}{m} \right) & \text{otherwise} \end{cases}$$

^۱ -Character-Based

^۲ -Longest Common Subsequence

$$d_w = \begin{cases} d_j & \text{if } d_j < b_t \\ d_j + (l_p(1-d_j)) & \text{otherwise} \end{cases}$$

m تعداد کل کاراکترهای مشترک و t تعداد کاراکترهایی است که جا به جا شده اند. در فرمول دوم که ضریب Jaro-Winkler است، d_j ضریب Jaro و ℓ تعداد کاراکترهای ثابت قبل اولین کاراکتر جابه جا شده و p ضریب ثابت تفکیک است که باید بیشتر از ۰,۲۵ باشد.

▪ نیدلمن-رانچ^۱

یک همترازی برابری بر دو ترتیب متوالی (مانند A و B) انجام می دهد. این الگوریتم نمونه ای از برنامه نویسی پویاست. این معیار برای جملات و رشته های هم اندازه دو شباهت قابل توجه کاربردی است.

▪ اسمیت واترمن^۲

این الگوریتم، الگوریتم نیدلمن رانچ را بهینه می کند به این صورت که به جای در نظر گرفتن تمام توالی کاراکتری، این الگوریتم سعی می کند با در نظر گرفتن بخشهای مختلف با همه طولهای ممکن میزان شباهت را بهینه کند. این الگوریتم برای متونی که به نظر شبیه نیستند اما ممکن است قسمت های کوچکی از آنها به هم شبیه باشد مناسب است.

▪ n-gram

زیر دنباله ای n موردی از دو رشته یا متن را با هم مقایسه می کند. مثلاً 1-gram کلمه به کلمه جملات را در نظر می گیرد یا 2-gram دو کلمه با هم از هر رشته، و به همین ترتیب فاصله و شباهت n-gram از هر رشته را تا حداکثر n اندازه گیری می کند.

۲,۱,۵,۲. شباهت بر پایه توالی^۳

▪ فاصله بلوک یا فاصله منهن

این الگوریتم فاصله بین دو نقطه روی یک مسیر مشبک را محاسبه می کند. در شباهت بین دو جمله این فاصله با اندازه گیری مجموع تفاوت بین اجزاء مترادف نظیر به نظیر بدست می آید.

$$d_1(p, q) = \|p - q\|_1 = \sum_{i=1}^n |p_i - q_i|$$

▪ شباهت کسینوسی

۱ اندازه گیری فاصله داخلی دو بردار متقاطع در فضا با استفاده از اندازه گیری کسینوس زاویه بین دو بردار. در اندازه گیری شباهت بین دو جمله یا رشته هر جمله نماد یک بردار

^۱-Needleman-Wunsch

^۲-Smith- Waterman

^۳-term-based

است و ارزش هر یک، برای نمونه تعداد دفعات، تکرار یک توالی، کاراکترها یا یک کلمه خاص در جمله، در نظر گرفته می شود.

$$similarity = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$

▪ ضریب تاس^۱

این ضریب به صورت دو برابر عبارات مشترک در رشته های مورد مقایسه تقسیم بر تعداد کل عبارت های در نظر گرفته شده در هر دو رشته محاسبه می شود.

$$Qs = \frac{2C}{A+B} = \frac{2|A \cap B|}{|A|+|B|}$$

▪ فاصله اقلیدسی^۲

یا فاصله L2 که همان فاصله بین دو نقطه در مختصات دکارتی است. برای شباهت بین دو جمله این فاصله جذر مجموع تفاوت بین عبارات متناظر در دو رشته در نظر گرفته

$$d(p,q) = \sqrt{(p_1-q_1)^2 + (p_2-q_2)^2 + \dots + (p_n-q_n)^2} = \sqrt{\sum_{i=1}^n (p_i-q_i)^2}$$

می شود.

▪ ضریب تشابه جاکارد^۳

نسبت مجموع عبارات مشابه و مشترک در دو رشته به عبارات غیر مشترک

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

▪ ضریب تطبیق^۴

یک ضریب تشابه بسیار ساده. اندازه عبارات مشترک به کل عبارات.

▪ ضریب همپوشانی^۵

شبیه به ضریب تاس با این تفاوت که شباهت بین دو رشته را در حالی اندازه می گیرد که کلمه ای از یکی زیرمجموعه ای از دیگری باشد.

$$overlap = \frac{|X \cap Y|}{\min(|X|, |Y|)}$$

¹-Dice's Coefficient

²-Euclidean distance

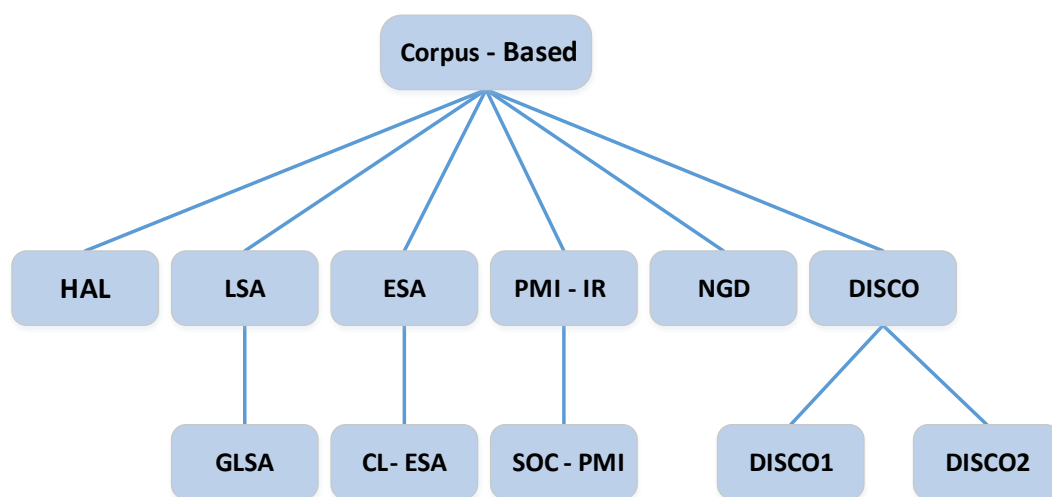
³-Jaccard similarity

⁴-Matching Coefficient

⁵-Overlap coefficient

۲,۵,۲. تشابه بر پایه مجموعه

سنجش تشابه دو متن بر پایه مجموعه یک معیار سنجش تشابه معنایی است، که شباهت بین کلمات را بر اساس اطلاعات مفید موجود در یک مجموعه بدست می آورد. این مجموعه حجم انبوهی از دست نوشته ها و متون گفتاری است که برای تحقیقات زبانی مورد استفاده قرار می گیرند. می توان گفت که میزان تشابه وابسته به مجموعه ای است که کلمات و جملات از آنها استخراج شده است. شکل ۲-۷ انواع روش های سنجش تشابه بر پایه مجموعه را نشان می دهد، که در ادامه به اختصار معرفی می شوند.



شکل ۲-۷- انواع الگوریتم های تشابه بر پایه مجموع

▪ (HAL) Hyperspace Analoue to Language

در این روش یک فضای معنایی از کلمات موجود در مجموعه ساخته می شود به این شکل که یک ماتریس کلمه در کلمه تشکیل می شود. به طوری که مقادیر آن نشان دهنده میزان ارتباط بین کلمه سطر و ستون متناظر با آن مقدار است. سپس آنتروپی پایین در ستون ماتریس توسط کاربر استخراج می شود. در تجزیه تحلیل متن برای استخراج ماتریس یک قالب n کلمه ای را در نظر می گیریم.

(مثلا $n=1$) سپس فاصله وقوع هر کلمه با کلمه اولی شمارش می شود. وزن مقادیر ماتریس رابطه معکوس با فاصله وقوع هر کلمه با کلمه ابتدایی دارد. کلمه ای که به کلمه ابتدایی نزدیک تر باشد از نظر معنایی رابطه نزدیکتری و در نتیجه وزن بیشتری دارد. HAL همچنین به ترتیب قرارگیری نیز توجه می کند. یعنی اینکه همسایه قبل از کلمه مورد نظر است یا بعد از آن به این منظور حالت همسایه مختلف را ثبت می کند.

▪ (LSA) Latent Semantic Analysis

تحلیل معنایی نهفته متداول ترین روش در میان روش های مبتنی بر مجموعه اصلی است. به ویژه برای تحلیل رابطه معنایی بین دستهای از متون مورد استفاده قرار می گیرد. روش کار به اختصار اینگونه است، که ابتدا ماتریس کلمه و سطر یا پاراگراف تشکیل می - شود، که ردیف ها نشان دهنده کلمات منحصر به فرد و ستون ها نشان دهنده سطرهاست. پس از روش تجزیه مقدار دهی منفرد (SVD)¹ به منظور کاهش سطر و ستون ماتریس در حالی که ساختار شباهت بین ردیف ها حفظ می شود، استفاده می کند. پس از آن کلمات با محاسبه کسینوس رابطه بین دو بردار تشکیل شده توسط هر دو ردیف مقایسه می شوند. مقادیر هر چه به یک نزدیک تر باشد، دو کلمه به هم شبیه تر بوده و هر چه به صفر نزدیک - تر باشد، شباهت کمتری دارند.

▪ تعمیم تحلیل معنایی نهفته GLSA²

یک چارچوب برای محاسبه میزان اهمیت عبارات و سند حاوی آن. در این روش که بر پایه LSA است، توجه بر ارزش سند حاوی عبارت به جای روش ارائه شده در LSA شده که به تعداد کلمات موجود در سطرها و مقایسه آنها توجه کرد. GLSA به این منظور نیاز به محاسبه میزان تشابه معنایی بین عبارات و نیز استفاده از روش های کاهش ابعاد مورد مقایسه دارد. برای رسیدن به نتیجه مطلوب بهترین و کارآمدترین روش ها را با هم ترکیب می کند و در آخر ماتریس عبارات موجود در اسناد که در LSA معرفی شده، برای محاسبه وزن در ترکیب بردارهای خطی از عبارات مورد استفاده قرار می گیرد.

▪ تحلیل معنایی صریح ESA³

معیاری برای محاسبه ارتباط معنایی بین دو متن دلخواه است. ESA به عنوان روشی برای بهبود طبقه بندی متون معرفی شده است. در این روش یک ماتریس tf-idf تشکیل می شود که مقادیر آن نشان دهنده رابطه معنایی میان دو متن است که توسط اندازه گیری کسینوسی محاسبه می شود.

¹-Singular Value Decomposition

²-Generalized Latent Semantic Analysis

³- Explicit Semantic Analysis

▪ تحلیل صریح معنایی متقابل زبانی CL-ESA¹

یک تعمیم چند زبانه از ESA است. که از یک متن چند زبانه به عنوان مرجع در مرکز چندین متن دیگر که از نظر زبانی مستقل هستند استفاده می کند. ارتباط این اسناد و میزان وابستگی آن ها به سن مرجع از طریق اندازه گیری کسینوسی محاسبه می شود. ویکیپدیا از این روش برای دسته بندی و اتصال اسناد متنی خود استفاده می کند.

▪ بازیابی اطلاعات متقابل نقطه به نقطه PMI-IR²

روشی است برای محاسبه شباهت بین کلمات است. موتور جستجو آستایستا از این روش برای محاسبه احتمال در جستجو پیشرفته استفاده می کند. دو کلمه که در یک صفحه وب نزدیک ترین موقعیت را نسبت به هم داشته باشد نمره تشابه بالاتری در PMI-IR خواهد داشت.

▪ بازیابی اطلاعات متقابل نقطه به نقطه نوع دوم SCO-PMI³

در این روش شباهت معنایی در واژه، با لیست کردن کلمات همسایه با آن در واژه در متن کلی انجام می شود. مزیت این روش در این است که می توان شباهت بین دو کلمه را که در همسایگی هم نیستند اما رابطه معنایی دارند را نیز اندازه گیری کرد.

▪ فاصله نرمال گوگل NGD⁴

اندازه گیری شباهت معنایی که از بازدیدهای صورت گرفته به وسیله موتور جستجو گوگل برای یک مجموعه از کلمات کلیدی بدست آمده است. در موتور جستجو فاصله دو کلمه کلیدی با معنی یکسان یا نزدیک به هم در زبان طبیعی رابطه نزدیکتری نسبت به دو کلمه غیر هم معنی دارند و در اصلاح به هم نزدیکتر هستند.

$$NGD(x, y) = \frac{\max\{\log f(x), \log f(y)\} - \log f(x, y)}{\log N - \min\{\log f(x), \log f(y)\}}$$

M در اینجا تعداد صفحات وب جستجو شده به وسیله Google، f(x) و f(y) تعداد موقعیت ها، در جستجو عبارتهای x و y هر دو وجود داشتند. اگر این دو کلمه به هم نزدیک نباشد و جدا از هم در صفحه ظاهر شوند، ضریب NGD بی نهایت است و چنانچه هر دو دائما کنار هم ظاهر شوند NGD آنها صفر یا برابر با مقداری بین مربع x و مربع y خواهد بود.

¹-Cross Language-Explicit Semantic Analysis

²-Pointwise Mutual Information-Information Retrieval

³-Pointwise Mutual InformationSecond Order Co-Occurrence

⁴-Normalized Google Distance

■ استخراج توزیعی کلمات مشابه با استفاده از تکرار وقوع کلمات DISCO¹

در این روش فرض بر این است که کلمات مشابه در زمینه های مشابه استفاده می شوند. بر همین اساس مجموعه های بزرگ متن مورد تجزیه و تحلیل آماری قرار می گیرند تا شباهت توزیعی بین کلمات استخراج شود. در DISCO شباهت توزیعی بین کلمات با استفاده از یک قاب متحرک با اندازه ± 3 برای اندازه گیری دفعات پدیدار شدن کلمه، اندازه گیری می شود. دو معیار DISCO1 و DISCO2 نیز توسعه ای از معیار اصلی هستند به طوری که DISCO1 شباهت دو کلمه را بر اساس ترتیب مجموعه که کلمه در آن است محاسبه می کند و DISCO2 میزان تشابه را بر اساس مجموعه کلمات توزیع شده مشابه با کلمات مورد نظر محاسبه می کند.

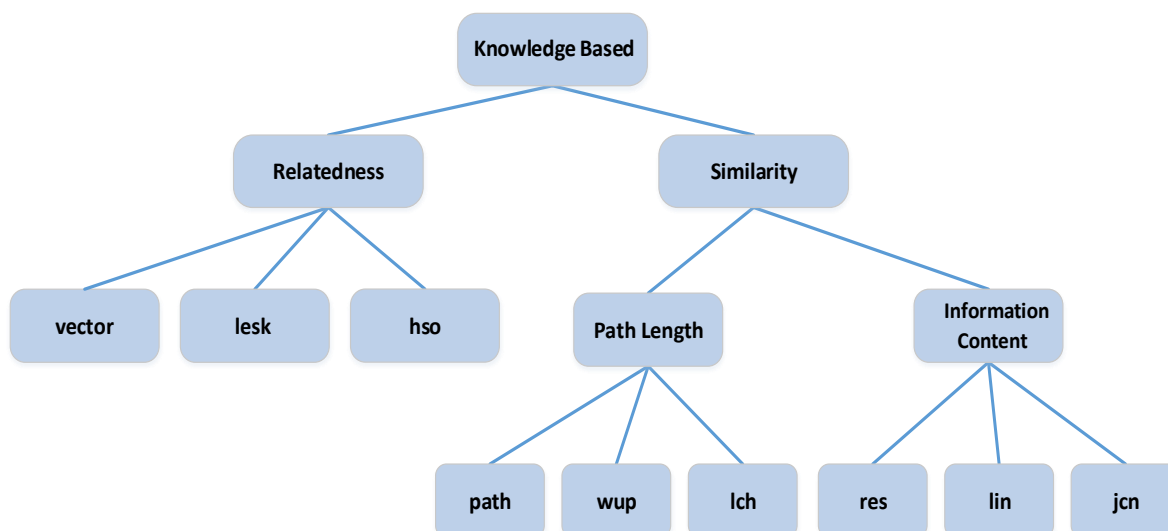
۳,۵,۲. تشابه بر پایه دانش

روش های مبتنی بر دانش بر اساس شناسایی میزان شباهت بین کلمات با استفاده از اطلاعات بدست آمده از شبکه های معنایی عمل می کنند. Wordnet² متداول ترین شبکه معنایی است. Wordnet یک پایگاه بزرگ واژه ای در زبان انگلیسی است. اسمها، فعل ها، صفت ها و قیدها به صورت مجموعه ای از مترادف های مشابه دسته بندی شده اند (synsets) که هر کدام بیان کننده یک مفهوم مجزا هستند.

synsets ها با استفاده از روابط معنایی و مفهومی و روابط لغوی به هم متصل هستند. روش های اندازه گیری تشابه بر پایه دانش به دو گروه تقسیم می شوند: اندازه گیری شباهت معنایی و اندازه گیری ارتباط معنایی. اندازه گیری شباهت معنایی همان گونه که از اسمش بر می آید هنگامی قابل اندازه گیری است، که دو کلمه در رابطه معنایی بر پایه شباهت با یکدیگر باشند، دو کلمه همانند و متشابه یکدیگر باشند اما اندازه گیری بر پایه رابطه معنایی دو کلمه می پردازد. به عنوان مثال اینکه کلمه ای نوع خاصی از دیگری باشد، یا دو کلمه مخالف هم باشند، یکی بخشی از دیگری باشد یا غیره.

¹ -Extracting DIStributionally similar words using CO-occurrences

² -www.wordnet.princeton.edu



شکل ۲-۸- انواع الگوریتم های تشابه برپایه دانش

از میان نه روش ارائه شده اندازه گیری شباهت معنایی (شکل ۲-۸) سه روش بر پایه محتوا اطلاعات و سه روش بر پایه طول مسیر اندازه گیری می شوند سه روش بر پایه رابطه مفهومی به اندازه گیری شباهت می پردازند. معیار path عددی را به عنوان میزان شباهت دو کلمه بر می گرداند که بر کوتاه ترین مسیر مفهومی که دو کلمه را به هم متصل می کند، دلالت دارد به عنوان مثال father و parent در شبکه جزء یک طبقه بندی متصل هستند یا نه و فاصله مفهومی بین این دو کلمه چقدر است. در روش HSO زنجیره ی لغوی بین دو کلمه پیدا شده و بر اساس آن رابطه دو کلمه در یکی از سه رشته ارتباطی تعریف می شود فوق العاده قوی- قوی- متوسط که حداکثر نتیجه برای یک رابطه عدد ۱۶ است. [۲۱]

۴,۵,۲. اندازه گیری شباهت ترکیبی

همانگونه که از اسم این روش ها بر می آید این روش ها با ترکیب روش های قبل سعی به از بین بردن نواقص و کاستی های هر روش با روش دیگر کرده اند تا معیارهای بهینه تری را ارائه دهند. بسیاری از تحقیق ها نیز در این حوزه صورت گرفته تا به حال هشت روش آزمایش شده ارائه شده که دوتای آنها بر پایه اندازه گیری بر اساس مجموعه، شش تای دیگر بر اساس اندازه گیری بر پایه دانش مطرح شده اند. روش ارائه شده در [۲۲] ابتدا شباهت معنایی بین کلمات از یک پایگاه دانش لغوی و مجموعه استخراج می شود و سپس در مرحله دوم تاثیر ترتیب و جای کلمه در معنای جمله را در نظر می گیرد.

در روش STS تشابه متن معنایی میزان شباهت بین کلمات را با ترکیب اطلاعات معنایی و نحوی اندازه گیری می کند. STS از دو روش شباهت خطی و شباهت معنایی به همراه روش انتخابی common word order بهره می گیرد.

STS در [۲۳] روی سی جفت کلمه روش جدید را آزمایش کرده و با محاسبه ضریب همبستگی پیرسون^۱ نتایج را بهبود بخشید. در روش [۲۴] نیز از اندازه گیری معنایی بر اساس مجموعه همراه با میزان شباهت معنایی بر اساس دانش برای کلمات هم نقش در جملات مختلف استفاده شده است. مهمترین خصوصیات این روش استفاده از مدل های یادگیری ماشین مانند رگرسیون خطی و مدل bugging برای بدست آوردن یک درجه شباهت موثر بین جملات است.

در [۲۵] دو روش اندازه گیری مفهومی با استفاده از wordnet و Ngram را با هم برای ارتباط دادن بین دو روش دستی و اتوماتیک انتخاب شده است.

همان طور که در قبل بیان شد، نمونه ها و بانک دادهای مورد استفاده در این تحقیق و به طور کل در مخازن خطا نرم افزار توسط یک گروه یا کاربر خاص تنظیم و ثبت نمی شوند. این داده ها توسط تمام کسانی که به نوعی در ارتباط با نرم افزار و پروژه هستند ثبت می شوند. پس طبیعی است که این متون از نظر نوع نگارش و دیکته لغات با الگو ویژه و اصولی نباشد. از سوی دیگر ممکن است این متون حاوی کلمات مشابه و هم معنا و حتی گاهی هم معنی اما غیر مربوط به هم باشند. در این شرایط برای استخراج دانش و فیلتر کردن داده ها نیاز به روشی است که این مشکلات تاثیر چندانی در نتایج آن نداشته باشد. همچنین در تشخیص جملات مشابه هم نیاز به دقت محاسباتی در روش های خطی و هم نیاز به دقت در معنا و نحو جملات و کلمات لازم است.

کلمات تخصصی در این متون پر اهمیت تر از کلمات متداول و اضافه هستند، پس نیاز است که اهمیت آنها به مراتب بیشتر از کلمات رایج در همه جملات است. روش مورد نظر پاسخگوی همه این نیاز ها خواهد بود همچنین در روش [۲۳] نه تنها شباهت ظاهری و معنایی و نحوی کلمات در نظر گرفته می شود بلکه به کلمات مجاور و حتی غیر مجاور آنها که در جمله یا متن ظاهر می شوند و روابط معنایی آنها با کلمه مورد نظر توجه می شود.

دلایل فوق باعث انتخاب روش اندازه گیری تشابه معنایی با استفاده از تشابه خطی و تشابه معنایی بر پایه مجموعه، در این تحقیق شده است. در ادامه به طور کامل این روش را توضیح می دهیم. روش مورد نظر شباهت بین دو متن را از نظر معنایی و اطلاعات نحوی (نظم متداول کلمات در زبان) مورد بررسی قرار می دهد. برای این کار از سه تابع شباهت استفاده می شود.

اول، شباهت رشته ها و شباهت معنایی کلمات محاسبه می شود. سپس برای ترکیب محاسبات با شباهت نحوی از تابع شباهت کلمات متداول انتخابی استفاده می کنیم. در نهایت ضریب شباهت دو متن با ترکیب شباهت رشته، شباهت معنایی بین کلمات و شباهت کلمات رایج و نرمال سازی محاسبات بدست می آید این روش STS^۲ نام دارد.

^۱ -Pearson correlation coefficient

^۲ -Semantic Text Similarity

فرض کنید کلمه ای در دو جمله به کار رفته باشد که این دو جمله معنا و حوزه مفهومی یکسانی داشته باشد اما در یک جمله کلمه مورد نظر اشتباه نوشته شده باشد.

به عنوان مثال دو جمله زیر را در نظر بگیرید

1.Einstein was a German-born theoretical physieist.

2.Einstain was the scientist of physics at 19 century.

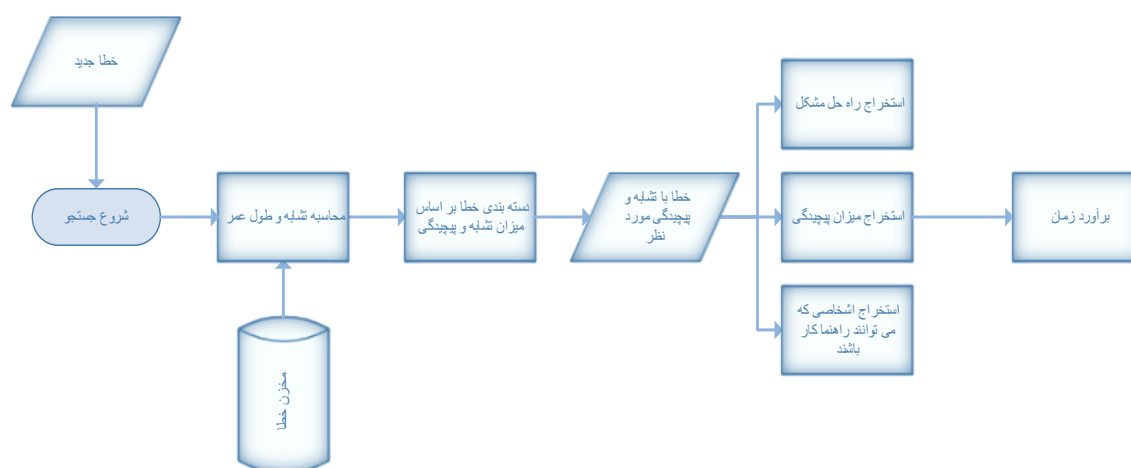
در دو جمله بالا Einstein در نقش اسم به کار رفته و هر دو جمله مفهومی نزدیک به هم دارند اما در جمله دوم این کلمه از نظر املائی درست نوشته نشده است. مگر از معیارهای تشابه مبنی بر فرهنگ لغات دیگر برای بررسی میزان تشابه دو جمله بالا استفاده کنیم به خاطر تفاوت دو کلمه، یا کلمات موجود در جمله میزان تشابه کمی را نشان می دهند در حالی که این دو جمله بر اساس رشته کلمات و مجموعه ای که کلمه مورد نظر در آن است شباهت را اندازه گیری کند. در ادامه به توضیح بخش های مختلف مورد استفاده در روش برگزیده خود می پردازیم.

استفاده از تکنیکهای داده کاوی برای کشف و رفع
خطاهای نرم افزار بر پایه ضریب تشابه معنایی متن
و خوشه بندی

۱.۳. مقدمه

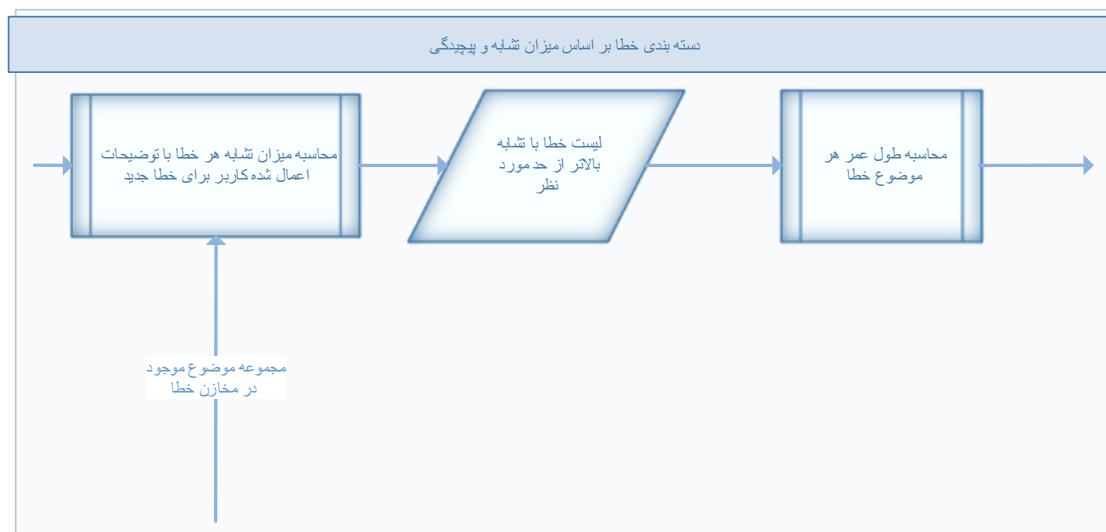
مخازن خطای نرم افزار سالهاست که به کمک مهندسين نرم افزار و مديران برای پياده سازی پروژه ها کوچک و بزرگ نرم افزاری ، متن باز يا بسته آمده است. هر خطا در اين مخازن به صورت فرآیندی جداگانه برای ثبت تمام رویدادها و اعمال نظرهای منطقی وموثر، درطول رفع مشکل و خطا از روی پروژه، از سوی افراد دخیل در پروژه ثبت می شود. این کار علاوه بر مزایای مربوط به پی-گیری روند رفع خطا ، مجموعه ای از دانش نهفته در این اطلاعات را در اختیار ذینفعان قرار می دهد. اگر این اطلاعات با حساسیت لازم و دقت کافی استخراج شود در روند تکمیل پروژه اثر به سزایی خواهد داشت. از آنجا که این داده ها بیشتر حاوی متون هستند نیاز به الگوریتم ها و معیارهایی برای بررسی متن احساس می شود. بیشتر این مخازن و کارهای قبلی انجام شده از روش های ساده و یا از موتورهای جستجوگر معمول برای استخراج دانش مفید استفاده کرده اند.

یکی از نکاتی که پایه این تحقیق و روش قرار گرفته است، تشابه بین متن توضیحی درباره خطای جدید و متون موجود در موضوع، توضیحات و دیگر متون ثبت شده مربوط به هر خطا است. این تشابه قابل اندازه گیری است و به یقین متن با تشابه بالاتر مربوط به خطا با تشابه بالاتر به خطای جدید است. این خطا شانس برطرف کردن خطای جدید را بالاتر برده و زمان رفع آن و در نتیجه هزینه های مربوط به آن را کاهش می دهد. به خاطر دلایلی که در فصل قبل بیان شد، انتخاب یک روش بر پایه تشابه معنایی ضروری است. در این فصل از روش محاسبه تشابه معنایی بر پایه مجموعه و تشابه خطی برای استفاده در یک بانک خطای نرم افزاری نمونه برای یافتن راه حل خطا و در نهایت محاسبه پیچیدگی خطا با استفاده از محاسبه طول عمر خطا استفاده می شود. شکل ۱-۳ مراحل مختلف روش ارائه شده در این تحقیق را نشان می دهد .



شکل ۱-۳- مراحل مختلف روش ارائه شده

در مرحله اول تعدادی خطا به عنوان نمونه کوچکی از یک بانک داده از داده های موجود در Bugzilla انتخاب شده. خلاصه، زمان ایجاد و زمان خاتمه (پاسخ قابل قبول) هر کدام استخراج می شود. در مرحله بعد میزان تشابه بین نمونه ها با جمله نمونه خطای جدید در پنج مرحله محاسبه و لیست این تشابهات که اعدادی بین ۰ و ۱ هستند استخراج می شود. در مرحله سوم طول عمر خطا ها محاسبه می شود این عدد پیچیدگی آنها را نشان می دهد. (شکل ۳-۲)



شکل ۳-۲- مراحل دسته بندی خطا

و در مرحله نهایی این دو لیست داده عددی خوشه بندی می شوند تا بتوان نمونه های منتخب را برای کمک به حل خطا و نیز پیش بینی مقادیر لازم مورد استفاده قرار می گیرند.

۲,۳. محاسبه تشابه معنایی بر پایه مجموعه و تشابه خطی

Islam در روش خود ۳ روش تشابه را برای برطرف کردن ضعف آنها ترکیب کرد. این روش روشی متداول در رفع نواقص الگوریتم هاست. برای محاسبه تشابه بین دو جمله یا متن در این روش سه مرحله کلی وجود دارد. ۱. محاسبه تشابه خطی بین کلمات ۲. محاسبه تشابه معنایی بر پایه مجموعه ۳. محاسبه تشابه بر اساس کلمات متداول ۴. محاسبه تشابه مجموع (محاسبه نهایی).

۱,۲,۳. شباهت خطی بین کلمات

یکی از روش های مورد انتظار در این نوع اندازه گیری روش طولانی ترین زیردنباله مشترک (LCS)^۱ با کمی تغییرات و نرمال سازی است. سه نسخه اصلاح شده متفاوت از LCS به کار رفته و سپس وزن های بدست آمده با هم جمع می شوند. به این نکته باید توجه کرد که نسخه های اصلاح شده LCS نتایج بهتری (دقت و جامعیت بالاتر) را نسبت به LCS و یا معیارهای شباهت دیگر خواهند داشت.

$$LCS(X_i, Y_j) = \begin{cases} 0 & \text{if } i=0 \text{ or } j=0 \\ LCS(X_{i-1}, Y_j) + 1 & \text{if } x_i = y_j \\ \max(LCS(X_i, Y_{j-1}), LCS(X_{i-1}, Y_j)) & \text{if } x_i \neq y_j \end{cases}$$

LCS – Length(X, Y)

```

1  m ← Length[X]
2  n ← Length[Y]
3  for i ← 1 to m
4      do c[i, 0] ← 0
5  for j ← 0 to n
6      do c[0, j] ← 0
7  for i ← 1 to m
8      do for j ← 1 to n
9          do if xi = yj
10             then c[i, j] ← c[i - 1, j - 1] + 1
11                 b[i, j] ← "□"
12             else if c[i - 1, j] ≥ c[i, j - 1]
13                 then c[i, j] ← c[i - 1, j]
14                     b[i, j] ← "↑"
15                 else c[i, j] ← c[i, j - 1]
16                     b[i, j] ← "←"
17  return c and b
```

¹ -Allison and Dix 1986

Melamed با تقسیم طولانی ترین زیر رشته مشترک به طول رشته طولانی دیگر سعی کرد که LCS را نرمال سازی کند. وی روش را LCSR نامید. اما در این روش طول رشته های کوتاهتر که گاهی اوقات در محاسبه شباهت دارای ارزش بالاتری بودند، در نظر گرفته نمی شود و این خود ضعف محسوب می شود.

اما در این روش LCS با در نظر گرفتن توامان طول رشته کوتاهتر و بلندتر نرمال سازی می شود (NLCS)^۱.

$$v_1 = NLCS(r_i, s_j) = \frac{length(LCS(r_i, s_j))^2}{length(r_i) \times length(s_j)}$$

در حالی که در LCS توالی کاراکترها مهم و ضروری نیست، اما این توالی برای داشتن درجه بالایی از تطبیق در رشته لازم است. این نکته در روش MCLS^۲ مورد توجه قرار گرفته است. در این روش ابتدا با توالی حداقل کاراکتر MCLS_۱ شروع و تا حداکثر کاراکتر MCLSn_۱ پیش می رود. الگوریتم ۱ MCLS_۱ در زیر نشان داده است. همچنین الگوریتم ۲ MCLSn_۱، روند کار را مشخص کرده است.

الگوریتم ۱: MCLS_۱ (LCS با حداکثر توالی n=1)

```

input :  $r_i, s_j, \lambda$  /* where  $|r_i| = \tau$ ,  $|s_j| = \eta$  and  $\tau \leq \eta$  */
output :  $r_i$ 
1  $\tau \leftarrow |r_i|$ ,  $\eta \leftarrow |s_j|$ 
2 while  $|r_i| \geq 0$  do
3   if  $r_i \cap s_j$  then /* i.e.,  $r_i \subset s_j = r_i$  */
4     return  $r_i$ 
5   else
6      $r_i \leftarrow r_i \setminus c_\tau$  /* i.e., remove the right most character from  $r_i$  */
7   end
8 end

```

¹-normalized longest common subsequence

²-maximal consecutive longest common subsequence starting at character

الگوریتم ۲: MCLSn (با حداکثر متوالی $n=n$)

```

input :  $r_i, s_j, \lambda$  /* where  $|r_i| = \tau$ ,  $|s_j| = \eta$  and  $\tau \leq \eta$  */
output :  $x$  /*  $x$  is the maximal consecutive LCS starting at any character  $n$  */
1  $\tau \leftarrow |r_i|$ ,  $\eta \leftarrow |s_j|$ 
2 while  $|r_i| \geq 0$  do
3 determine all  $n$ -grams from  $r_i$  where  $n = 1 \dots |r_i|$  and
4  $\bar{r}_i$  is the set of  $n$ -grams
5 if  $x \in s_j$  where  $\{x | x \in \bar{r}_i, x = \max(\bar{r}_i)\}$  then
6 return  $x$ 
7 else
8  $r_i \leftarrow r_i \setminus x$  /* remove  $x$  from  $\bar{r}_i$  */
9 end
10 end

```

در الگوریتم اول دو رشته به عنوان ورودی گرفته حداکثر رشته متوالی مشترک در هر دو رشته که با کاراکتر اول رشته کوتاه تر آغاز شده باشد به عنوان خروجی برگردانده می شود. اگر الگوریتم ۲ همین کار را برای تک تک کاراکترها تکرار کرده و بزرگترین رشته متوالی از هر دو رشته را به عنوان خروجی باز می گرداند. این دو الگوریتم نیز نرمالسازی می شوند.

$$v_2 = NMCLCS(r_i, s_j) = \frac{\text{length}(MCLCS(r_i, s_j))^2}{\text{length}(r_i) \times \text{length}(s_j)} \quad \text{فرمول ۱-۳}$$

$$v_3 = NMCLCS_n(r_i, s_j) = \frac{\text{length}(MCLCS_n(r_i, s_j))^2}{\text{length}(r_i) \times \text{length}(s_j)} \quad \text{فرمول ۲-۳}$$

مقادیر وزنی بدست آمده از سه قبلی v_1 و v_2 و v_3 با ضرایب وزنی w_1 و w_2 و w_3 که $w_1 + w_2 + w_3 = 1$ با هم جمع شده و میزان شباهت دو جمله بدست می آید.

$$\alpha = w_1 v_1 + w_2 v_2 + w_3 v_3$$

در این جا برای w ها ارزش برابر در نظر گرفته می شود. به عنوان مثال فرض کنید که روش را روی دو رشته زیر پیاده کنیم.

$ri = fcmfkzyw$

$sj = fcfmfkzny$

$LCS(ri, sj) = fcmfkz$

$MCLCS1(ri, sj) = fc$

$MCLCSn(ri, sj) = mfzkz$

$NLCS(ri, sj) = 7^2 / (8 \times 9) = 0.68$

$NMCLCS1(ri, sj) = 2^2 / (8 \times 9) = 0.056$

$NMCLCSn(ri, sj) = 4^2 / (8 \times 9) = 0.22$

$\alpha = w_1v_1 + w_2v_2 + w_3v_3 = 0.33 \times 0.68 + 0.33 \times 0.056 + 0.33 \times 0.22 = 0.32$ میزان شباهت

۲،۲،۳. تشابه معنایی بین کلمات

در روش مورد بررسی توجه بر روش های اندازه گیری بر پایه مجموعه است، زیرا این روش ها انواع گسترده ای از متون را پوشش می دهند متونی که در دنیای واقعی استفاده می شود در حالی که در روش های بر پایه دانش به این نکته توجه نشده است.

PMI-IR یک روش ساده برای محاسبه شباهت متنی مجموعه ای از کلمات است که با استفاده

از اطلاعات نقطه به نقطه متقابل به صورت زیر تعریف می شود

$$PMI(w_1, w_2) = \log p(w_1 \text{ AND } w_2) / p(w_1)p(w_2)$$

w_1 و w_2 کلمات و $p(w_1 \text{ AND } w_2)$ احتمال رخ داد توامان دو واژه است. اگر دو کلمه از نظر آماری

مستقل باشند $p(w_1 \text{ AND } w_2) = p(w_1).p(w_2)$ و در غیر این صورت امکان پدیدار شدن همزمان را با هم

در متن دارند و $p(w_1 \text{ AND } w_2) > p(w_1).p(w_2)$.

PMI-IR توسط موتور جستجوگر آلتاویستا^۱ برای جستجو نحوی و محاسبه احتمالات استفاده

می شود. در ساده ترین حالت دو کلمه هنگامی که در یک سند یافت شوند، احتمال رخ داد همزمان

را دارند. این احتمال با بازیابی اسناد قابل اندازه گیری به صورت تقریبی است.

$$PMI - IR(w_1, w_2) = hits(w_1 \text{ AND } w_2) / hits(w_1) hits(w_2)$$

$hits(x)$ تعداد متونی است که حاوی کلمه x بوده اند.

روش رخ داد همزمان مرتبه دوم (SOC-PMI)^۲ لیستی از کلمات همسایه با دو کلمه مورد

نظر در یک مجموعه بزرگ را، بر اساس اطلاعات متقابل نقطه به نقطه مرتب می کند.

مزیت دیگر این روش این است که شباهت میان دو کلمه در هر جای متن نیز می توان اندازه

گیری کرد. حتی اگر در کنار هم نباشد. تنها کافی است کلمات همسایه مشابه داشته باشند.

^۱ -www.altavista.com

^۲ - Second Order Co-occurrence PMI (SOC-PMI)

مزیت دیگر این روش (PMI-IR) این است که شباهت میان دو کلمه در هر جای متن نیز می توان اندازه گیری کرد. حتی اگر در کنار هم نباشند. تنها کافی است که کلمات همسایه مشابه داشته باشند.

در این روش از مجموعه ملی بریتانیا (BNC) به عنوان یک منبع برای متون استفاده می شود. کلماتی که در هر دو فهرست مشترک است لیست شده و جمع ارزش PMI (از لیست مقابل) برای محاسبه برای بدست آوردن نسبت ارزش معنایی محاسبه می شود. تابع نقطه به نقطه اطلاعات متقابل برای کلماتی که دارند، تعریف می شود.

$$f^{pmi}(t_i, w) = \log_2 \frac{f^b(t_i, w) \times m}{f^t(t_i) f^t(w)} \quad \text{فرمول ۳-۳}$$

$f^t(t_i)$ تعداد دفعاتی است که t_i در مجموعه ظاهر می شود. $f^b(t_i, w) > 0$ تعداد دفعاتی را که t_i و w در یک قاب^۱ ظاهر می شود را نشان می دهد. m تعداد کل توکن ها در مجموعه است. برای کلمه w مجموعه کلمات X^w به صورت زیر تعریف می شود.

$$f^{pmi}(t_1, w) \geq f^{pmi}(t_2, w) \geq \dots \geq f^{pmi}(t_{\beta-1}, w) \geq f^{pmi}(t_{\beta}, w)$$

قاعده کلی برای β به صورت زیر است

$$\beta = \left(\log(f^t(w)) \right)^{\frac{2 \log_2(n)}{\mu}} \quad \text{فرمول ۳-۴}$$

μ یک مقدار ثابت است که در این روش ۶/۵ در نظر گرفته شود. μ با اندازه مجموعه رابطه مستقیم دارد. هر چه μ کوچکتر باشد کلمات مهمتر ممکن است نادیده گرفته شود.

تابع مجموع PMI، β کلمه w_1 نسبت به کلمات دیگر به صورت زیر تعریف می شود.

$$f(w_1, w_2, \beta) = \sum_{i=1}^{\beta} (f^{pmi}(X_i^{w_1}, w_2))^{\gamma} \quad \text{فرمول ۳-۵}$$

اگر $f^{pmi}(X_i^{w_1}, w_2) > 0$ تمام مقادیر مثبت PMI برای کلمات مجموعه X^{w_2} در X^{w_1} نیز هستند. به عبارت دیگر کلماتی که از نظر معنایی به w_2 نزدیک باشند. مطمئناً در لیست کلمات هم معنی

$$\text{sim}(w_1, w_2) = \frac{f^{\beta}(w_1)}{\beta_1} + \frac{f^{\beta}(w_2)}{\beta_2} = \frac{44.255}{23} + \frac{22.364}{19} = 3.10$$

مقدار γ باید مقداری بزرگ تر از یک باشد. پس تابع مجموع β -PMI برای w_1 نسبت به w_2 با

$\beta = \beta_1$ برابر و برای w_2 نسبت به w_1 با $\beta = \beta_2$ برابر ۳. در نهایت تابع شباهت معنایی PMI بین دو کلمه w_1 و w_2 به صورت زیر تعریف می شود [۲۳].

^۱ در اینجا قاب ۱۱ کلمه فرض شده است
^۲ هر چه این مقدار بیشتر باشد توجه به کلمات با ارزش بیشتر است. در آزمایشهای متوالی در این روش به این نتیجه رسیده اند که $\gamma = 3$ مقدار مناسب است. $\gamma \geq 4$ میزان تاکید بیهوده و تاخیر زیاد را به همراه داشت و کلمات کم

$$\text{sim}(w_1, w_2) = \frac{f(w_1, w_2, \beta_1)}{\beta_1} + \frac{f(w_1, w_2, \beta_2)}{\beta_2} \quad \text{فرمول ۳-۶}$$

الگوریتم ۳ نحوه نرمال سازی مقادیر بدست آمده از $\text{sim}()$ را نشان می دهد.

الگوریتم ۳. نرمال سازی محاسبات شباهت معنایی

input : set D, distance measure dist, number k of cluster

output : A partitioning P of the set D of documents

- 1 choose randomly k data points from D as starting centroids $t_{p_1}^- \dots t_{p_k}^-$
- 2 repeat
- 3 Assign each point of P to the closest centroid with respect to dist
- 4 (Re-) calculate the cluster centroids $t_{p_1}^- \dots t_{p_k}^-$ of clusters $P_1 \dots P_k$
- 5 until cluster centroids $t_{p_1}^- \dots t_{p_k}^-$ are stable
- 6 return set $\mathbf{P} := \{P_1 \dots P_k\}$ of clusters.

مثال زیر روند کار را توضیح می دهد.

با استفاده از روش توضیح داده در بالا میزان تشابه معنایی دو کلمه car و automobile در ۱۲ جمله جدول ۳-۱ محاسبه می شود.

جدول ۳-۱-۱۲ جمله به عنوان نمونه داده

1	pursuit accident claim car driver exclude
2	soak motorist company car driver risky
3	company car driver tend travel farther
4	job engineer disappear fall mechanical engineer car industry worst affect
5	Sign recession car industry
6	brightest engineer moment car industry
7	yugoslavia benefit direct investment automobile industry
8	acreage expand emergence automobile industry
9	automobile industry among hardest hit recession
10	automobile industry largely male force
11	component supplier automobile industry expand

ارزش تر نادیده گرفته می شود Islam و Inkpen در بخش کوچکی از BNC پس از تکرار آزمایش تعداد مناسب β ، μ و γ را بدست آوردند.

$f^i(t_i)$ تعداد تکرار کلمات اصلی محاسبه می شود. (جدول ۳-۲)

جدول ۳-۲ - تعداد تکرار کلمات اصلی

t_i	$f^i(t_i)$	t_i	$f^i(t_i)$
disappear	1	worst	1
yugoslavia	1	soak	1
Pursuit	1	fall	1
brightest	1	supplier	1
travel	1	company	2
Benefit	1	recession	2
risky	1	farther	1
Sign	1	car	6
male	1	investment	1
accident	1	industry	10
affect	1	force	1
mechanical	1	job	1
claim	1	client	1
among	1	Tend	1
moment	1	hardest	1
engineer	3	component	2
automobile	6	manufacturer	1
emergence	1	expand	2
direct	1	driver	3
hit	1	exclude	1
Largely	1		

در مرحله بعد $f^b(t_i, w)$ برای همه t_i ها برای دو کلمه مورد نظر محاسبه می شود. قاب با طول ۱۱

کلمه (± 5) برای این در نظر گرفته شده است. (جدول ۳-۳)

جدول ۳-۳ - محاسبه $f^b(t_i, w)$ برای همه کلمات جدول ۲-۳

t_i	$f^b(t_i, W_1)$	t_i	$f^b(t_i, W_2)$
brightest	1	acreage	1
accident	1	Emergence	1
affect	1	direct	1
motorist	1	hit	1
disappear	1	largely	1
worst	1	yugoslavia	1
pursuit	1	supplier	1
soak	1	benefit	1
fall	1	recession	1
travel	1	male	1
risky	1	Investment	1
recession	1	industry	7
company	2	force	1
sign	1	client	1
farther	1	hardest	1
mechanical	1	Component	2
claim	1	expand	2
tend	1	manufacturer	1
industry	3	Among	1
moment	1		
engineer	3		
exclude	1		
driver	3		

$$f^b('engineer', w_2) = 7$$

$$f^b('engineer', w_1) = 3$$

تابع اطلاعات متقابل نقطه به نقطه (PMI) برای t_i ها یی که $0 < f^b(t_i, w)$ بر اساس فرمول ۳-۳ محاسبه می شود. (جدول ۳-۴)

جدول ۳-۴ - مقدار تابع اطلاعات متقابل نقطه به نقطه برای تمامی مقادیر با $0 < f^b(t_i, w)$

X_i (also t_i)	$f^{pmi}(t_i, W_1)$	Y_i (also t_i)	$f^{pmi}(t_i, W_2)$
motorist	3.544	emergence	3.544
disappear	3.544	direct	3.544
worst	3.544	acreage	3.544
pursuit	3.544	hit	3.544
soak	3.544	largely	3.544
travel	3.544	yugoslavia	3.544
brightest	3.544	supplier	3.544
fall	3.544	benefit	3.544
risky	3.544	male	3.544
company	3.544	investment	3.544
sign	3.544	among	3.544
farther	3.544	force	3.544
accident	3.544	client	3.544
affect	3.544	hardest	3.544
mechanical	3.544	component	3.544
tend	3.544	manufacturer	3.544
claim	3.544	expand	3.544
engineer	3.544	industry	3.029
moment	3.544	recession	2.544
driver	3.544		
exclude	3.544		
recession	2.544		
industry	1.807		

β_1 و β_2 طبق فرمول ۳-۳ برای محاسب نهایی محاسبه می شود.

$$\beta_1 = 24.88 \quad \beta_2 = 24.88$$

μ در اینجا ۰,۷ در نظر گرفته شده است. γ نیز برای محاسبات همان مقدار بهینه ۳ در نظر گرفته می شود. دو کلمه recession و industry هر دو جز مجموعه Y_i و X_i هستند پس طبق فرمول ۳-۴ و ۵-۳

$$f^\beta(W_1) = (f^{pmi}(\text{"recession"}, W_2))^\gamma + (f^{pmi}(\text{"industry"}, W_2))^\gamma = (2.544)^3 + (3.029)^3 = 44.255$$

$$f^\beta(W_2) = (f^{pmi}(\text{"industry"}, W_1))^\gamma + (f^{pmi}(\text{"recession"}, W_1))^\gamma = (1.807)^3 + (2.544)^3 = 22.364$$

$$\text{sim}(w_1, w_2) = \frac{f^\beta(w_1)}{\beta_1} + \frac{f^\beta(w_2)}{\beta_2} = \frac{44.255}{23} + \frac{22.364}{19} = 3.10$$

۳,۲,۳. تشابه جملات بر اساس عبارات مشترک

اگر دو متن کلمات مشابه با هم داشته باشند می توان میزان مشابهت آنها را بر اساس این کلمات مشترک بدست آورد (این کلمات یا همگی در یک موقعیت مکانی هستند، و یا در موقعیت های مختلف در جملات ظاهر می شوند). هستینگز [۲۶] معتقد بود که این کلمات ارزش چندانی در اندازه گیری شباهت معنایی جملات کوتاه ندارد. در این روش برای آنکه از اهمیت کم این معیار چشم پوشی نشود فاکتور وزنی آن را کمتر از ۰/۵ در نظر می گیریم .

روش کار به این صورت است که ، فرض کنیم P و R دو جمله با کلمات مشابه هستند و $|R| < |P|$ ، تمام کلمات مشابه در دو جمله استخراج می شوند. اگر X مجموع کلمات مشابه در P و Y مجموع کلمات مشابه در R باشد هر کدام از این مجموعه ها ترتیب مخصوص به خود دارند. از آنجا که تعداد کلمات P بیشتر از R است و در اصطلاح P بزرگتر از R است. به کلمات موجود در مجموعه X به ترتیب موقعیتشان در جمله وزنی از ۱ تا δ اختصاص می دهیم و همین وزن به کلمات مجموعه Y الحاق می شود. سپس میزان شباهت این جمله بر اساس دستور زیر محاسبه می شود.

$$X = \{x_1, x_2, \dots, x_\delta\}$$

$$Y = \{y_1, y_2, \dots, y_\delta\}$$

$$S_o = \frac{1 - |x_1 - y_1| + |x_2 - y_2| + \dots + |x_\delta - y_\delta|}{|1 - x_\delta| + |x_2 - x_{\delta-1}| + \dots + |x_\delta - x_1|}$$

$$S_o = \begin{cases} 1 - \frac{2 \sum_{i=1}^{\delta} |x_i - y_i|}{\delta^2} & \text{if } \delta \text{ is even} \\ 1 - \frac{2 \sum_{i=1}^{\delta} |x_i - y_i|}{\delta^2 - 1} & \text{if } \delta \text{ is odd and } \delta > 1 \\ 1 & \text{if } \delta \text{ is odd and } \delta = 1 \end{cases}$$

فرمول ۳-۷

به عنوان مثال دو جمله را در نظر بگیرید

P: Einstein was a German-born theoretical physicist

R: The theoretical physicist – Einstein lived at 19 century

$X = \{\text{Einstein, theoretical, physicist}\}$ $X = \{1, 2, 3\}$

$Y = \{\text{theoretical, physicist, Einstein}\}$ $Y = \{2, 1, 3\}$

$$S_o = 1 - \frac{2 \times 2}{9} = \frac{5}{9} = 0.55$$

با توجه به فرمول ۷-۳

۴,۲,۳. شباهت کلی جملات

Islam روش خود را در ۶ مرحله ارائه خلاصه کرده است:

۱- ابتدا تمام کلمات اضافه در جملات برای پیدا کردن کلمات کلیدی پاک می شود اگر P و

R دو جمله مورد نظر باشند m کلمه از P و n کلمه از R کلمات اصلی ما را تشکیل

می دهند. (حذف حروف و کلمات اضافه)

$$R = \{r_1, r_2, \dots, r_\delta\} \quad P = \{p_1, p_2, \dots, p_\delta\} \quad n \geq m$$

۲- در این مرحله کلمات مشابه در این مجموعه علامت گذاری می شوند. δ کلمه مشابه در R

و P کنار گذاشته شده و بقیه برای بررسی نگه داشته می شود. اگر $m = \delta$ بود به مرحله ۶

میرویم. در غیر این صورت $\delta < m$ به ترتیب ادامه می دهیم.

$$R = \{r_1, r_2, \dots, r_{n-\delta}\} \quad P = \{p_1, p_2, \dots, p_{m-\delta}\}$$

۳- ماتریس تشابه خطی دو مجموعه $a(m-\delta) \times (n-\delta)$ به صورت زیر تشکیل می شود. هر

α_{ij} در a به صورت زیر محاسبه می شود:

اگر $p_i \in P$ کاراکتر τ داشته باشد $p_i = \{c_1, c_2, \dots, c_\tau\}$ و هر $r_j \in R$ ، η کاراکتر

$r_j = \{r_1, r_2, \dots, r_\eta\}$ و $\tau \leq \eta$ به صورتی که τ طول کوتاه ترین کلمه و η بلندترین کلمه مشترک

است [۲۳].

$$v_1 \leftarrow \text{NLCS}(p_i, r_j)$$

$$v_3 \leftarrow \text{NMCLCSn}(p_i, r_j)$$

$$v_2 \leftarrow \text{NMCLCSI}(p_i, r_j)$$

$$\alpha_{ij} \leftarrow w_1 v_1 + w_2 v_2 + w_3 v_3 \quad \text{فرمول ۳-}$$

۸

$$w^1 + w^2 + w^3 = 1$$

$$M_1 = \begin{pmatrix} \alpha_{11} & \alpha_{12} & \dots & \alpha_{1j} & \dots & \alpha_{1(n-\delta)} \\ \alpha_{21} & \alpha_{22} & \dots & \alpha_{2j} & \dots & \alpha_{2(n-\delta)} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \alpha_{i1} & \alpha_{i2} & \dots & \alpha_{ij} & \dots & \alpha_{i(n-\delta)} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \alpha_{(m-\delta)1} & \alpha_{(m-\delta)2} & \dots & \alpha_{(m-\delta)j} & \dots & \alpha_{(m-\delta)(n-\delta)} \end{pmatrix}$$

مثال زیر این مرحله را توضیح می دهد

$$r_j = \text{"allmileage_make_maxkm"}$$

$$p_i = \text{"make_minmile_distance_possible_take"}$$

اجرا مرحله اول

$$r_j = \{\text{all, mileage, make, max, km}\} \quad m=5$$

$$p_i = \{\text{make, min, mile, distance, possible, take}\} \quad n=6$$

اجرا مرحله دوم

$$r_j = \{\text{all, mileage, max, km}\}$$

$$p_i = \{\text{min, mile, distance, possible, take}\}$$

	Min	mile	distance	possible	take
all	0	0.055	0.041	0.027	0.082
mileage	0.188	0.565	0.058	0.076	0.058
max	0.11	0.082	0.027	0	0.055
km	0.11	0.082	0	0	0.123

به عنوان نمونه میزان شباهت بین possible و mileage اینگونه محاسبه شده:

$$v_1 = 3^2 / (8 \times 7) = 0.16$$

$$v_2 = 0$$

$$v_3 = 2^2 / (8 \times 7) = 0.071$$

$$\alpha_{24} = 0.33 \times v_1 + 0.33 \times v_2 + 0.33 \times v_3 = 0.076$$

مطابق فرمول ۳-۸

۴- ماتریس $M_2 = (\beta_{ij})_{(m-\delta) \times (n-\delta)}$ ، ماتریس تشابه معنایی از روی الگوریتم ۳ تشکیل می

شود.

$$M_2 = \begin{pmatrix} \beta_{11} & \beta_{12} & \dots & \beta_{1j} & \dots & \beta_{1(n-\delta)} \\ \beta_{21} & \beta_{22} & \dots & \beta_{2j} & \dots & \beta_{2(n-\delta)} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \beta_{i1} & \beta_{i2} & \dots & \beta_{ij} & \dots & \beta_{i(n-\delta)} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \beta_{(m-\delta)1} & \beta_{(m-\delta)2} & \dots & \beta_{(m-\delta)j} & \dots & \beta_{(m-\delta)(n-\delta)} \end{pmatrix}$$

۵- ماتریس M از طریق دو ماتریس بالابه صورت زیرتشکیل می شود. $M \leftarrow \psi M_1 + \phi M_2$.
 ψ فاکتور وزنی شباهت خطی و ϕ فاکتور وزنی شباهت معنایی است. که $\phi + \psi = 1$. هر
 کدام از این ضرایب اگر صفر در نظر گرفته شوند ماتریس مربوطه در محاسبه ماتریس
 مجموع M لحاظ نمی شود. اگر این ضرایب هر دو ۰,۵ باشند یعنی شباهت خطی و معنایی
 هر دو به یک میزان اهمیت دارند. عنصر با ماکزیمال ارزش γ_{ij} ، در ماتریس M استخراج می
 شود.

اگر این مقدار بزرگتر از صفر بود به لیست ρ اضافه می شود، $\rho \leftarrow \gamma_{ij} \cup \rho$. سطر و ستون
 مربوط به آن از ماتریس حذف می شود. این کار تا زمانی که $\gamma_{ij} = 0$ و یا $m - \delta - |\rho| = 0$
 ادامه پیدامی کند.

۶- تمام مقادیر P و $(1 - w_f + w_f S_o)$ هم جمع می شوند. طبق فرمول $S(P, R)$ میزان
 شباهت چند گانه با ضرب در مجموع m و n و تقسیم بر دو برابر ضریب آنها به مقداری
 بین صفر و یک بالانس می شوند. در این روش می توان از هر کدام از روش ها و مقادیر
 شباهت صرفه نظر کرد. مثلاً از معیار شباهت کلمات متداول با صفر اعلام کردن S_o صرفه
 نظر کرد و یا اینکه با صفر در نظر گرفتن w_f اهمیت نحوی را نادیده گرفت [۲۳].

۵,۲,۳. نمونه توضیح روند کار

دو تکه متن مورد نظر برای مقایسه و محاسبه ضریب شباهت به روش فوق با R و P مشخص شده
 است.

$$R = \text{"allmileage_make_maxkm"}$$

$$P = \text{"make_minmile_distance_possible_take"}$$

مرحله ۱: جدا کردن کلمات اصلی و مشخص

$$m=5$$

$$R = \{all, mileage, make, max, km\}$$

$$n=6$$

$$P = \{make, min, mile, distance, possible, take\}$$

مرحله ۲: کلمه هایی که در هر دو مجموعه مشترک است حذف می شوند.

$$R = \{all, mileage, make, max, km\}$$

$$P = \{make, min, mile, distance, possible, take\}$$

مرحله ۳: ماتریس تشابه خطی با محاسبه تشابه خطی تک تک کلمات بر اساس فرمول ۳-۸

تشکیل می شود.

$$M_1 = \begin{matrix} & \begin{matrix} Min & mile & distance & possible & take \end{matrix} \\ \begin{matrix} all \\ mileage \\ max \\ km \end{matrix} & \begin{bmatrix} 0 & 0.055 & 0.041 & 0.027 & 0.082 \\ 0.188 & 0.565 & 0.058 & 0.076 & 0.058 \\ 0.11 & 0.082 & 0.027 & 0 & 0.055 \\ 0.11 & 0.082 & 0 & 0 & 0.123 \end{bmatrix} \end{matrix}$$

مرحله ۴: ماتریس تشابه معنایی نیز تشکیل می شود.

$$M_2 = \begin{matrix} & \begin{matrix} Min & mile & distance & possible & take \end{matrix} \\ \begin{matrix} all \\ mileage \\ max \\ km \end{matrix} & \begin{bmatrix} 0.172 & 0.233 & 0.48 & 0 & 0.813 \\ 0.587 & 0.976 & 0.826 & 0 & 0.558 \\ 0.199 & 0.194 & 0.141 & 0 & 0.243 \\ 0.67 & 0.962 & 0.89 & 0 & 0.408 \end{bmatrix} \end{matrix}$$

مرحله ۵: ماتریس M از طریق دو ماتریس قبل به کمک فرمول زیر تشکیل می شود.

$$M \leftarrow \psi M_1 + \phi M_2$$

$$M = \begin{matrix} & \begin{matrix} Min & mile & distance & possible & take \end{matrix} \\ \begin{matrix} all \\ mileage \\ max \\ km \end{matrix} & \begin{bmatrix} 0.086 & 0.144 & 0.26 & 0.013 & 0.447 \\ 0.388 & 0.771 & 0.442 & 0.038 & 0.308 \\ 0.154 & 0.138 & 0.084 & 0 & 0.194 \\ 0.39 & 0.522 & 0.445 & 0 & 0.266 \end{bmatrix} \end{matrix}$$

سپس مقدار ماکزیمم به مجموعه ρ اضافه می شود و سطر و ستون آن حذف می شود. این کار

تا زمانی که $\gamma_{ij} = 0$ و یا $m - \delta - |\rho| = 0$ ادامه پیدا می کند.

$$M = \begin{matrix} & \begin{matrix} Min & mile & distance & possible & take \end{matrix} \\ \begin{matrix} all \\ mileage \\ max \\ km \end{matrix} & \begin{bmatrix} 0.086 & 0.144 & 0.26 & 0.013 & 0.447 \\ 0.388 & 0.771 & 0.442 & 0.038 & 0.308 \\ 0.154 & 0.138 & 0.084 & 0 & 0.194 \\ 0.39 & 0.522 & 0.445 & 0 & 0.266 \end{bmatrix} \end{matrix}$$

$$\rho = \{0.771\}$$

$$M = \begin{matrix} & \begin{matrix} Min & distance & possible & take \end{matrix} \\ \begin{matrix} all \\ max \\ km \end{matrix} & \begin{bmatrix} 0.086 & 0.26 & 0.013 & 0.447 \\ 0.154 & 0.084 & 0 & 0.194 \\ 0.39 & 0.445 & 0 & 0.266 \end{bmatrix} \end{matrix}$$

$$\rho = \{0.771, 0.447\}$$

$$M = \begin{matrix} & \begin{matrix} Min & distance & possible \end{matrix} \\ \begin{matrix} max \\ km \end{matrix} & \begin{bmatrix} 0.154 & 0.084 & 0 \\ 0.39 & 0.445 & 0 \end{bmatrix} \end{matrix}$$

$$\rho = \{0.771, 0.447, 0.445\}$$

$$\begin{aligned} & \text{Min distance} \\ M &= \max [0.154 \quad 0] \\ \rho &= \{0.771, 0.447, 0.445, 0.154\} \end{aligned}$$

مرحله ۶: و در نهایت محاسبه ضریب تشابه ترکیبی برای دو جمله R و S

$$\begin{aligned} \text{SimilarityScore}(R_i, S_j) &= \frac{(\delta + \sum_{i=1}^{|p|} \rho_i) \times (m+n)}{2mn} && \text{فرمول ۳-۹} \\ &= (1 + 1.817) \times 11/60 = 0.516 \end{aligned}$$

۳,۳. خوشه بندی داده ها

۱,۵,۲. الگوریتم K-means

با توجه به [۲۷] الگوریتم K-means یک الگوریتم ساده تکرار شونده است، که مجموعه داده ها را به K خوشه تقسیم می کند. این الگوریتم در سالهای مختلف در زمینه های گوناگون توسط افرادی مانند Lloyd (۱۹۵۷، ۱۹۸۲)، Forgey (۱۹۶۵)، Friedman و Rubin (۱۹۶۷) و McQueen (۱۹۶۷) توسعه یافته است. الگوریتم روی یک مجموعه از بردار های d بعدی تعریف می شود، که $D = \{\mathbf{x}_i \mid i = 1, \dots, N\}$ و $\mathbf{x}_i \in \mathbb{R}^d$ نشان دهند i امین نقطه داده است. تکنیک مورد استفاده برای انتخاب نمونه های اولیه برای تشکیل K دسته، شامل نمونه گیری تصادفی از مجموعه داده هاست. سپس الگوریتم دو قدم زیر را به صورت مکرر انجام می دهد.

▪ گام اول، تخصیص داده:

هر نقطه داده ها به نزدیک ترین مرکز جرم اختصاص داده می شود. اختصاص داده با استفاده از یک سری روابط قراردادی و دلخواه انجام می گیرد. این فرآیند منجر به دسته بندی داده ها می شود.

▪ گام دوم، جابجایی دسته ها:

نمایندگی در دسته به مرکزی ترین نقطه در دسته اختصاص داده می شود. اگر میزان احتمال داده ها یک اندازه باشد. پس جابجایی ها به اندازه انتظار یعنی به اندازه میانگین وزنی خواهد بود. با توجه به این مطلب الگوریتم اگر انتساب داده ها ارزش Cj در معادله تغییر زیادی نداشته باشد، همگرا خواهد بود.

$$\sum_{i=1}^n \left(\underset{j}{\operatorname{argmin}} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2 \right)$$

اصول الگوریتم در روند زیر دنبال شده است.

الگوریتم ۴ : K-means

input : set D , distance measure $dist$, number k of cluster

output : A partitioning P of the set D of documents

1 choose randomly k data points from D as starting centroids $t_{p_1}^- \dots t_{p_k}^-$

2 repeat

3 Assign each point of P to the closest centroid with respect to $dist$

4 (Re-) calculate the cluster centroids $t_{p_1}^- \dots t_{p_k}^-$ of clusters $P_1 \dots P_k$

5 until cluster centroids $t_{p_1}^- \dots t_{p_k}^-$ are stable

6 return set $P := \{P_1 \dots P_k\}$ of clusters.

یک مسئله مهم در استفاده از الگوریتم، مسئله تعیین و اندازه گیری کمیت میزان نزدیکی، در مرحله واگذاری داده به دسته ها است. به طور پیش فرض اندازه گیری میزان نزدیکی، فاصله اقلیدسی در نظر گرفته می شود. الگوریتم معمولاً بسیار حساس به مرکز جرم است. انتخاب مرکز جرم مجدد برای گروه، مشکل واگرایی را برای داده های پراکنده حل می کند.

منابع و مآخذ

၁. Naresh Kumar Nagwani, Pradeep Singh/"Bug Mining Model Based on Event-Component Similarity to Discover Similar and Duplicate GUI Bugs"/ 2009 IEEE International Advance Computing Conference (IACC 2009)/Patiala, India/6-7 March 2009.
၂. Bergroth, L. Dept. of Comput. Sci., Turku Univ., Finland ; Hakonen, H. ; Raita, T/"A survey of longest common subsequence algorithms"/String Processing and Information Retrieval, Proceedings/ Seventh International Symposium on/ SPIRE 2000.
၃. Islam, A. and Inkpen, D. /" Second Order Co-occurrence PMI for Determining the Semantic Similarity of Words"/ in Proceedings of the International Conference on Language Resources and Evaluation (LREC)/ Genoa, Italy/ pp. 1033–1038/ 2006.
၄. Naresh Kumar Nagwani, Shrish Verma/ "Predicting Expert Developers for Newly Reported Bugs Using Frequent Terms Similarities of Bug Attributes"/ 2011 Ninth International Conference on ICT and Knowledge Engineering / 2011 IEEE.
၅. Hui Zeng, David Rine/ " Estimation of Software Defects Fix Effort Using Neural Networks"/ Proceedings of the 28th Annual International Computer Software and Applications Conference (COMPSAC'04)/ 2004 IEEE.
၆. Junzo Watada/ "Analysis of Software Reliability by Fuzzy Regression Model"/ TENCON 2000
၇. Lucas D. Panjer/ " Predicting Eclipse Bug Lifetimes"/ Fourth International Workshop on Mining Software Repositories (MSR'07)/ 2007 IEEE.
၈. Suma.V, Pushpavathi T.P, and Ramaswamy.V/"An Approach to Predict Software Project Success by Data Mining Clustering"/ International Conference on Data Mining and Computer Engineering (ICDMCE'2012)/ Bangkok,(Thailand)/ December 21-22, 2012.
၉. Cathrin Weiß ,Thomas Zimmermann, Rahul Premraj , Andreas Zeller ,Saarland University/ " How Long will it Take to Fix This Bug?"/ Fourth International Workshop on Mining Software Repositories (MSR'07)/ 0-7695-2950-X/07 \$20.00 © /2007 IEEE
၁၀. Naresh Kumar Nagwani, Shrish Verma/" Predictive Data Mining Model for Software Bug Estimation Using Average Weighted Similarity"/ Advance Computing Conference (IACC)/ 2010 IEEE 2nd International.
၁၁. Naresh Kumar Nagwani, Ashok Bhansali/"A Data Mining Model to Predict Software Bug Complexity Using Bug Estimation and Clustering"/ 2010 International Conference on Recent Trends in Information/ Telecommunication and Computing/ 2010 IEEE.
၁၂. Naresh Kumar Nagwani, Shrish Verma/ "Predicting Expert Developers for Newly Reported Bugs Using Frequent Terms Similarities of Bug Attributes"/ 2011 Ninth International Conference on ICT and Knowledge Engineering / 2011 IEEE.

۱۳. B. Berliner/ "Parallelizing software development In Proceedings of the USENIX"/ Winter 1990 Technical Conference/ volume 341, page 352, 1990
۱۴. C. M. Pilato, B. Collins-Sussman, and B. W. Fitzpatrick/ "Version Control with Subversion"/ O'Reilly Media/ 2008
۱۵. www.sourceforge.net
۱۶. www.code.google.com
۱۷. Stephen W. Thomas/ "Mining Software Repositories with Topic Models"/ School of Computing Queen's University Kingston, Ontario, Canada. Technical Report 2012-586. IEEE.
۱۸. Golnoosh Abaee, Department of Studies in Computer Science, Islamic Azad University, Roodehen Branch, Iran. D.S.Guru, Department of Studies in Computer Science, University of Mysore, Manasagangothri, Mysore, 570006, India/ "Enhancement of Bug Tracking Tools; the Debugger"/ 2nd International Conference on Software Technology and Engineering (ICSTE)/ 2010.
۱۹. Bug life cycle / www.softwaretestinghelp.com/bug-life-cycle/
۲۰. Naresh Kumar Nagwani, Dr. Shrish Verma 1 Assistant Professor, Computer Science & Engineering; 2 Associate Professor & Head, Electronics & Tel. Communication Engg. Mm National Institute of Technology Raipur, 1nknagwani.cs@nitrr.ac.in, 2shrishverma@nitrr.ac.in/ On Studying the Effect of Sample Size in Evaluation of Bug Classifiers/ ISSN: 0974-6846, Vol: 6 Issue: 1/ January 2013
۲۱. Wael H. Gomaa, Computer Science Department Modern Academy for Computer Science & Management Technology Cairo, Egypt. Aly A. Fahmy, Computer Science Department Faculty of Computers and Information, Cairo University Cairo, Egypt/ "A Survey of Text Similarity Approaches"/ International Journal of Computer Applications/ (0975 – 8887) Volume 68– No.13/ April 2013.
۲۲. Li, Y., McLean, D., Bandar, Z., O'Shea, J., & Crockett, K/ "Sentence similarity based on semantic nets and corpus statistics"/ IEEE Transactions on Knowledge and Data Engineering/ 18(8), 1138–1149 / 2006.
۲۳. Islam, A., & Inkpen, D/ "Semantic text similarity using corpus-based word similarity and string similarity"/ ACM Transactions on Knowledge Discovery from Data/ 2(2), 1–25/ 2008.
۲۴. Nitish, A., Kartik, A. & Paul, B. DERI&UPM/ "Pushing Corpus Based Relatedness to Similarity: Shared Task System Description"/ First Joint Conference on Lexical and Computational Semantics (*SEM)/ pages 643–647/ Montreal, Canada, Association for Computational Linguistics/ June 7-8, 2012.

٢٥. Davide, B., Ronan, T. Nathalie A.&Josiane, M. (2012), IRIT/ "Textual Similarity Combining Conceptual Similarity with an N-Gram Comparison Method"/ First Joint Conference on Lexical and Computational Semantics (*SEM)/ pages 552–556/ Montreal, Canada/June 7-8, 2012 Association for Computational Linguistics.

٢٦. Naresh Kumar Nagwani, Dr. Shrish Verma¹Assistant Professor, Computer Science & Engineering; ²Associate Professor & Head, Electronics & Tel. Communication Engg. Mm National Institute of Technology Raipur, ¹knagwani.cs@nitrr.ac.in, ²shrishverma@nitrr.ac.in/ On Studying the Effect of Sample Size in Evaluation of Bug Classifiers/ ISSN: 0974-6846, Vol: 6 Issue: 1/ January 2013

٢٧. XindongWu , Vipin Kumar , J. Ross Quinlan , Joydeep Ghosh , Qiang Yang , Hiroshi Motoda , Geoffrey J .McLachlan , Angus Ng , Bing Liu , Philip S. Yu , Zhi-Hua Zhou , Michael Steinbach , David J.Hand , Dan Steinberg / "Top 10 algorithms in data mining/ Springer-Verlag London Limited "/ Published online: 4 December 2007.

٢٨. Osama Abu Abbas ,Computer Science Departmentm/"Comparisons Between Data Clustering Algorithms"/The International Arab Journal of Information Technology/ Vol. 5, No. 3,/July 2008.

٢٩. Rui Xu, Student Member, IEEE and Donald Wunsch II/"Survey of Clustering Algorithms"/ IEEE Transactions on neural networks, VOL. 16, NO. 3, May 2005