

طراحی سامانه تشخیص دستبرد ادبی جمله‌بنیاد در متون فارسی به کمک هم‌جوشی گواها

حمید آهنگر بهان و غلامعلی منتظر

دانشکده مهندسی صنایع و سیستم‌ها، دانشگاه تربیت مدرس، تهران، ایران

چکیده

امروزه، می‌توان با رونوشت‌برداری از منابع وسیع اینترنتی سند جدیدی را تهیه و به نام خود ارائه داد که مصداقی از دستبرد ادبی است. حفظ حقوق مالکیت مادی و معنوی نویسندگان ضرورت تشخیص اسناد اصلی و اسناد رونوشت‌برداری شده را ایجاب می‌کند. برای این منظور از الگوریتم‌های تشخیص دستبرد ادبی و سامانه‌هایی که مبتنی بر این الگوریتم‌ها طراحی شده‌اند استفاده می‌شود. عملکرد سامانه‌های تشخیص دستبرد ادبی که در واقع میزان شباهت بین دو متن را ارزیابی می‌کنند، بر اساس دو شاخص دقت یافته‌های مشابه و مدت زمان تشخیص سنجیده می‌شود. تاکنون سنج‌های مختلفی برای ارزیابی شباهت دو سند ارائه شده که کارایی آنها به محتوای متن و منابع مورد استفاده مانند واژه‌نامه برای مقایسه بین واژه‌های دو سند محدود است. در این مقاله روشی ارائه شده است که از نظریه گواها برای هم‌جوشی اطلاعات، به منظور ارزیابی تشابه دو سند و کشف دستبرد ادبی با توجه به کیفی و ناکامل بودن عوامل اثرگذار بر سنجش شباهت بین دو متن، استفاده می‌کند. روش پیشنهادی در گام نخست جمله‌های موجود در سند را به دو بخش عمومی و تخصصی تقسیم کرده و سپس با استفاده از سنج‌های متفاوت و استفاده از «هستان‌نگار تخصصی» امتیاز تشابه برای هر بخش را محاسبه و در نهایت در دو سطح، میزان شباهت بین دو سند را استنتاج می‌کند؛ در سطح نخست نتایج سنج‌های شباهت‌سنجی به‌عنوان گواها (با باور پایه مشخص) با قاعده دمپستر-شفر با هم ترکیب شده و به‌عنوان گواهی جدید به سطح دوم منتقل می‌شوند. در سطح دوم نتیجه سطح نخست و گواها جدید از طریق قاعده میانگین‌گیری ترکیب شده و توابع باور و مقبولیت نهایی محاسبه و شباهت بین دو جمله (سند) ارزیابی می‌شود. سامانه مذکور بر روی بیکره فارسی که با استفاده از مقاله‌های کنفرانس سالانه یادگیری و آموزش الکترونیکی تولید شده، مورد ارزیابی قرار گرفته است که با دقت بیش از ۸۶٪ امکان شناسایی اسناد مشابه را دارد.

واژگان کلیدی: دستبرد ادبی، هم‌جوشی داده، سنج‌های شباهت‌سنجی، نظریه گواها، شباهت‌سنجی معنایی.

۱- مقدمه

با فراگیر شدن شبکه وب و حجم زیاد اطلاعاتی که این شبکه در اختیار قرار می‌دهد، چالش‌های بسیاری در حوزه پردازش زبان‌های طبیعی به وجود آمده است. یکی از روش‌های مؤثر برای کاهش و سازماندهی این حجم زیاد داده‌ها، شناسایی اسناد مشابه و دسته‌بندی آنهاست (ونگ و هوجز، ۲۰۰۶). با یافتن اسناد مشابه امکان هر نوع استفاده غیرقانونی، که دستبرد ادبی نمونه بارز آن است، محدود خواهد شد.

در تشخیص اسناد جعلی از اسناد اصل، انتخاب درست ویژگی‌ها و روش‌های تشابه‌یابی متن جنبه کلیدی دارد (هرمان و زاکا، ۲۰۰۶). به همین دلیل نیاز است که

فرآیند به‌کارگیری اسناد توسط کاربران به‌دقت بررسی شود. در حال حاضر، با توسعه شبکه‌های اطلاعاتی و به‌ویژه وب، کاربران به راحتی می‌توانند سند جدیدی را با رونوشت‌برداری از منابع موجود به نام خود تهیه کنند که طبیعتاً مصداقی از دستبرد ادبی محسوب می‌شود. البته گاهی کاربر برای رد گم‌کردن از واژه‌های مترادف یا تغییر ساختار جمله‌ها استفاده می‌کنند که در این صورت تشخیص این نوع دستبرد به‌صورت خودکار مشکل است و سامانه‌هایی همانند «گپس»^۱ (عثمان و همکاران، ۲۰۱۳)، «اسکم»^۲ (رام و همکاران، ۲۰۱۴) و «چک»^۳ (کایا و همکاران، ۲۰۱۴)، قابلیت

¹ COPS

² SCAM

³ CHECK

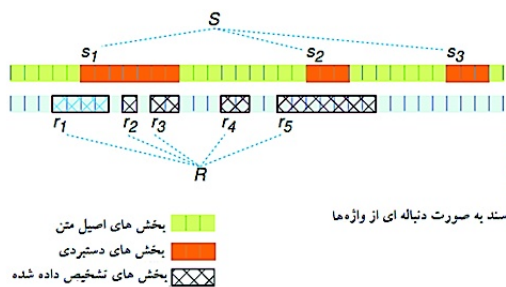
ادامه این مقاله بدین صورت تنظیم شده است: در بخش ۲ مسئله دستبرد ادبی تعریف می‌شود. در بخش ۳ مروری بر کارهای پیشین در زمینه سامانه‌های شناسایی دستبرد ادبی، روش‌های شباهت‌سنجی و در بخش ۴ مفاهیم اصلی نظریه گواه ارائه و روش‌های ترکیب گواه در این نظریه بحث می‌شود. در بخش ۵ معماری پیشنهادی برای سامانه شباهت‌سنجی بیان می‌شود؛ پیاده‌سازی این سامانه و نتایج عملی حاصل از کاربرد آن در بخش ۶ بیان خواهد شد و در نهایت بخش ۷ به نتیجه‌گیری می‌پردازد.

۲- تعریف دستبرد ادبی

دستبرد ادبی^۴ در ساده‌ترین شکل به صورت «رونوشت‌برداری از اسناد بدون تأیید مرجع اصلی» تعریف می‌شود که یکی از شایع‌ترین موارد آن، استفاده از مطالب دیگران در نوشته‌ای به نام خود است (جوی، ۱۹۹۹).

فرآیند تشخیص دستبرد ادبی را به کمک نمادهای ریاضی به صورت زیر می‌توان تعریف کرد (استین و همکاران، ۲۰۰۹): مجموعه مرجع D و مجموعه اسناد مشکوک D_q مفروض است، فرآیند تشخیص دستبرد ادبی عبارت است از یافتن بخش‌های مشکوک s_q از D_q که مشابه بخش‌های s_x در مجموعه مرجع D است.

شکل (۱) این فرآیند را نشان می‌دهد. در این فرآیند سعی بر این است تا بخش‌های تشخیص‌داده‌شده بیش‌ترین تطابق را با بخش‌های دستبردی داشته باشد. به‌عنوان نمونه در شکل ۱، بخش ۳ دستبردی تشخیص داده نشده است:



شکل-۱: فرآیند تشخیص دستبرد ادبی

انواع مختلف دستبرد ادبی عبارتند از (موشکه و گیپ، ۲۰۱۳):
 ۱- دستبرد تحت اللفظی^۵: این نوع دستبرد عبارت است از رونوشت برداری متن بدون تغییر یا با تغییرات کم.

شناسایی آن را ندارند. اشکال اصلی این سامانه‌ها، کندی، وابستگی به ساختار لغوی و عدم توجه به ساختار معنایی جمله‌ها و واژه‌هاست (عثمان و همکاران، ۲۰۱۳؛ باو و همکاران، ۲۰۰۳). بنابراین تشخیص متون هم‌معنائوسی^۱ شده، برای آنها بسیار دشوار خواهد بود. به همین دلیل چالش اصلی در این زمینه یافتن الگوریتمی مناسب برای بهبود دقت یافته‌های مشابه در مقایسه دو متن است که ساختار لغوی و معنایی جمله را همزمان در نظر داشته باشد. شایان توجه است که اگرچه مدت زمان سنجش شباهت دو متن نیز می‌تواند از شاخص‌های ارزیابی سامانه‌های تشخیص دستبرد ادبی باشد، اما از آنجا که مسئله دستبرد ادبی به‌طور معمول به صورت برون‌خط انجام می‌شود، دقت تشخیص بیش از مدت زمان آن مورد توجه است. سامانه‌های دستبرد ادبی از روش‌های اندازه‌گیری شباهت بین اجزای دو سند، چه در سطح لغوی و چه در سطح معنایی، استفاده می‌کنند از طرفی در دنیای واقعی، هر سند مجموعه‌ای از محتواهای متنوع و مختلف است و باید روشی را یافت که اطلاعات محتوای متن را به‌خوبی نشان دهد و سنجش شباهت‌سنجی و منبع‌اطلاعاتی مورد نیاز در همان زمینه و محتوا را در روش‌های شباهت‌سنجی به کار برد. از آنجا که محتوای سند به چندین بخش (از جمله تخصصی و عمومی) تقسیم می‌شود، به همین دلیل لازم است، سامانه، از منابع اطلاعاتی مختلف برای تصمیم‌گیری درباره شباهت بین دو سند استفاده کند تا دیدی جامع‌تر و معقول‌تر نسبت به مسئله داشته و در نهایت اطلاعات حاصل از هر بخش را با هم ترکیب کند. هم‌جوشی داده‌ها راه‌حلی برای یک‌پارچه‌سازی اطلاعات حاصل از منابع مختلف است و می‌توان آن را به صورت «فرآیند چندسطحی و چندوجهی برای تشخیص، جمع‌بندی، اصلاح، برآورد و ترکیب خودکار داده و اطلاعات از منابع مختلف» تعریف کرد (وایست، ۱۹۹۱؛ دونگ و همکاران، ۲۰۰۹). این روش در مسائل بسیاری، مانند پردازش تصویر، پردازش سیگنال و استنتاج مورد استفاده قرار گرفته است (لنز و همکاران، ۲۰۱۲).

در این مقاله پس از بررسی سامانه‌های تشخیص دستبرد ادبی موجود، مزایا و معایب آنها مطرح می‌شود و سپس روشی مبتنی بر هم‌جوشی داده‌ها^۲ و به‌کمک نظریه گواه^۳ و در قالب سامانه تشخیص دستبرد ادبی برای شناسایی اسناد مشابه فارسی ارائه خواهد شد.

¹ Paraphrase
² Data Fusion
³ Evidence Theory

⁴ Plagiarism
⁵ Literal plagiarism

(معنایی و یا لغوی) و نیز دقت در اندازه‌گیری شباهت مشخص می‌شود. به‌طورمعمول مقدار این آستانه در سامانه‌های تشخیص تشابه متغیر در نظر گرفته شده و بر اساس خواست کاربر تنظیم می‌شود. یکی از سامانه‌های طراحی‌شده برای شناسایی شباهت بین متون، سامانه «کُپس» است که در آن توابع رمزگذاری^۹ برای تولید اثرانگشت^{۱۰} به کار گرفته شده و مقادیر رمزنگاری دو سند برای تشخیص شباهت استفاده می‌شود (برین و همکاران، ۱۹۹۵). سامانه «اسکم» نیز از روش‌های بازیابی اطلاعات مبتنی بر واژه برای تشخیص رونوشت‌برداری استفاده می‌کند که تنها قابلیت شناسایی اسناد کوتاه را دارد (شیواکومار و همکاران، ۱۹۹۵). سامانه «چک» با هدف کاهش پیچیدگی، ابتدا تعداد کمی از ویژگی‌های دو سند را مورد بررسی قرار می‌دهد و در صورت برقرار بودن شرایط نخستین الگوریتم روی کل دو سند اجرا می‌شود (لنژ و همکاران، ۱۹۹۷). در سامانه «تشخیص رونوشت‌برداری از کالاهای دیجیتال»^{۱۱} نیز از روش‌های مبتنی بر بردار، برای تشخیص شباهت استفاده شده است (باو و همکاران، ۲۰۰۱). برای بررسی جمله به جمله دو سند نیز سامانه «اسپلت»^{۱۲} پیشنهاد شده که با محاسبه تعداد تکرار واژه‌ها این کار را انجام می‌دهد (کالبرگ و همکاران، ۲۰۰۳). روش‌های دیگری از جمله، «نشان‌گذاری دیجیتالی»^{۱۳} که اطلاعات اضافه‌ای را به سند ضمیمه می‌کند نیز برای استفاده در این سامانه‌ها پیشنهاد شده است (هیاری و همکاران، ۲۰۰۵؛ زینی و همکاران، ۲۰۰۶). روش‌های خوشه‌بندی اسناد نیز یکی از روش‌های بازیابی اطلاعات است که بازیابی داده را با کاهش زمان جستجو در جایابی اسناد، خلاصه‌سازی متون و کاهش زمان مقایسه در تشخیص دستبرد ادبی بهبود می‌بخشد. در برخی پژوهش‌ها نیز از واژه‌های کلیدی برای یافتن خوشه‌های مشابه بین اسناد استفاده شده است (برین و همکاران، ۱۹۹۵).

پژوهش‌های بسیاری نیز برای تشخیص دستبرد ادبی بر ویژگی‌های ساختاری متن مانند سرسندها^{۱۴}، بخش‌ها، پاراگراف‌ها و مراجع تمرکز دارند. چاو و رحمان از روش‌های نمایش متن به صورت ساختار درختی استفاده کرده و به کمک شبکه نگاشت خودسازمانده چندلایه^{۱۵} آن را غنی و

«رونوشت برداری و چسباندن»^۱ از متون منابع دیگر، اصلی‌ترین شکل این نوع دستبرد است (و بر- وولف، ۲۰۱۱).

۲- دستبرد پنهانی^۲: در این نوع دستبرد از شیوه‌هایی برای پنهان کردن رونوشت برداری متن استفاده می‌شود. سه نوع دستبرد پنهانی وجود دارد:

۱-۲- قطعه‌کردن و چسباندن^۳: در این شیوه از رونوشت برداری و تجمیع تکه متن‌ها با تغییرات جزئی (مانند تغییر ترتیب واژه‌ها، جایگزینی آنها با مترادف‌شان و یا حذف و اضافه کردن واژه‌ها) استفاده می‌شود (و بر- وولف، ۲۰۱۱).

۲-۲- هم‌معنائویسی^۴: عبارت است از بازنویسی خودآگاهانه ایده‌های دیگران به صورت لغوی و نگارش آن بدون اعلام منبع اصلی (لنکستر، ۲۰۰۳).

۳-۲- پنهان‌کاری فنی^۵: در این روش متونی ساخته می‌شود که برای ماشین قابل تشخیص نباشد. در این شیوه از روش‌هایی مانند جایگزینی نویسه‌ها با نمادهای گرافیکی و الفبای خارجی یا اضافه کردن تصادفی کلمات با قلم سفید استفاده می‌شود (هیتر، ۲۰۱۰).

۳- دستبرد ترجمه‌ای^۶: این نوع دستبرد ادبی عبارت است از تبدیل خودکار یا دستی و آگاهانه متن از یک زبان به زبان دیگر بدون بیان منبع اصلی آن (و بر- وولف، ۲۰۱۱).

۴- دستبرد ایده^۷: شامل استفاده از مفاهیم و ایده‌های دیگران مانند تصاحب دیدگاه‌های پژوهشی بدون اطلاع‌رسانی درباره منبع آن (مورر و همکاران، ۲۰۰۶).

۵- خوددستبرد^۸: به معنای استفاده کامل یا قسمتی از نوشته‌های خود بدون ارجاع به اصل (براتگ و محمود، ۲۰۰۹).

۳- روش‌های شناسایی دستبرد ادبی

در این بخش کارهای پژوهش‌گران در زمینه دستبرد ادبی و شباهت‌سنجی متون بررسی می‌شود. دو سند را مشابه گوییم اگر میزان شباهت واژه‌های آنها چه از نظر لغوی و چه معنایی از آستانه‌ای معلوم بیشتر باشد (گراوند و همکاران، ۲۰۱۴). این آستانه با توجه به نوع رویکرد سنجش

⁹ Hashing

¹⁰ Fingerprint

¹¹ Copying Detection System of Digital Goods

¹² SPLaT

¹³ Digital Watermarking

¹⁴ Header

¹⁵ Multi-Layer Self-Organizing Maps Model

سال ۱۳۹۵ شماره ۱ پیاپی ۲۷

¹ Copy and paste

² Disguised plagiarism

³ Shake and paste

⁴ Paraphrasing

⁵ Technical disguise

⁶ Translated plagiarism

⁷ Idea plagiarism

⁸ Self-plagiarism

شباهت متفاوتی را برای زوج‌واژه‌ها در نظر بگیرند، به همین دلیل اگر سندی شامل واژه‌های یک حوزه خاص در اختیار داشته باشیم، این سنج‌ها به‌تنهایی کارایی مناسبی نخواهند داشت؛ لذا در این‌گونه متون استفاده از سنج‌های ترکیبی برای رفع این معایب احساس می‌شود. جدول (۲) بیان‌گر دسته‌بندی سامانه‌های پیشنهادی برای تشخیص دستبرد ادبی است. این جدول نشان می‌دهد تاکنون تمرکز پژوهش‌گران به استفاده از روش‌های شباهت‌سنجی لغوی بوده و کمتر از شباهت‌سنجی معنایی استفاده کرده‌اند. سنج‌های شباهت‌سنجی معنایی نیز به‌شدت وابسته به ساختار و کیفیت پایگاه دانش است. به همین دلیل در این مقاله به‌دنبال طراحی سامانه‌ای برای تشخیص دستبرد ادبی هستیم که بتواند مشکلات موجود در روش‌های شباهت‌سنجی لغوی و معنایی را پوشش داده و از مزایای هر دو روش نیز استفاده کند.

۴- نظریه گواه دمپستر - شفر

شناخت عوامل اثرگذار بر مسائل مهندسی در دنیای امروزی، اصلی‌ترین عنصر حل آن خواهد بود. در دنیای واقعی هر داده و اطلاعاتی به‌طور معمول با نوعی نقصان^۴ روبه‌رو است که برای حل آنها باید از روش‌های متناسب با نوع نقص آن استفاده کرد. در حال حاضر نظریه‌های زیادی برای نمایش نقص داده مانند نظریه احتمال^۵، نظریه مجموعه فازی^۶، نظریه امکان^۷، نظریه مجموعه نادقیق^۸ و نظریه گواه دمپستر شفر مورد استفاده قرار می‌گیرد. بسیاری از این نظریه‌ها توانایی بازنمایی داده ناقص را تنها در منطری خاص دارند. برای مثال توزیع احتمال، قابلیت بیان عدم اطمینان در داده و نظریه مجموعه فازی توانایی بازنمایی داده مبهم با استفاده از تابع عضویت فازی را دارد. این نظریه تنها توانایی حل مسائل بالبهام را به‌وسیله متغیرهای زبانی و هم‌جوشی از طریق قواعد فازی دارد درحالی‌که نظریه گواه، توانایی هم‌جوشی داده نامطمئن و نامعلوم را دارد و در حل انواع دیگر داده ناقص ناتوان و برای هم‌جوشی داده با تعارض بالا نیز ناکارآمد است (خالقی و همکاران، ۲۰۱۳).

برای تشخیص دستبرد ادبی به کار برده‌اند (چاو و رحمان، ۲۰۰۹). در روش پیشنهادی اقبال و همکاران، از رویکرد مرسوم مبتنی بر بردار و همچنین روش‌های مبتنی بر گراف برای تشخیص پاراگراف‌های مشکوک استفاده شده است (اقبال و همکاران، ۲۰۱۲). گراوند و همکاران نیز روشی مبتنی بر فیلترهای بلوم^۱ را برای بررسی دستبرد ادبی در اطلاعات شخصی و محرمانه برای حفظ اطلاعات افراد پیشنهاد کرده‌اند (گراوند، ۲۰۱۴).

اصلی‌ترین بخش هر سامانه شناسایی دستبرد ادبی، الگوریتم شباهت‌سنجی آنهاست از این‌رو در این بخش مزایا و معایب هر رویکرد برای درک بهتر مسئله تشخیص دستبرد ادبی، بیان می‌شود:

تشابه متن را می‌توان در دو حوزه کلی «تشابه لغوی^۲» و «تشابه معنایی^۳» بررسی کرد. تشابه‌یابی لغوی، به ظاهر واژه‌ها توجه می‌کند و برای سنجش آن در ابتدا واژه‌ها ریشه‌یابی شده و سپس از طرق مختلفی تشابه اندازه‌گیری می‌شود (متزلر و همکاران، ۲۰۰۷). «تشابه معنایی» نیز مفهومی است که تشابه بین مجموعه‌ای از اسناد یا واژه‌ها را از طریق معنای آنها مشخص می‌کند (می‌هالی‌شا و همکاران، ۲۰۰۶). رویکردها و روش‌های تشابه‌یابی و همچنین مقایسه این روش‌ها با شاخص‌هایی که در پیاده‌سازی نرم‌افزاری آنها و همچنین تشخیص دستبرد ادبی مهم است، در جدول (۱) نمایش داده شده است. ضعف اصلی روش‌های تشابه‌یابی لغوی این است که نمی‌توانند عبارات متفاوت را، که از نظر معنایی با هم مرتبط هستند، شناسایی کنند؛ در تشابه‌یابی معنایی نیز چالش اصلی پایگاه دانش و پیکره غنی و جامع مورد نیاز است؛ به‌عنوان نمونه سنج‌های مبتنی بر پایگاه دانش، فرآیند رمزنگاری خاصی را برای روابط معنایی دارند، در مقابل، سنج‌های مبتنی بر پیکره، چون از داده‌های خام و بی‌ساختار استفاده می‌کنند به‌صورت گسترده قابل استفاده هستند؛ ولی متأسفانه این ویژگی‌ها برای این دو دسته سنج‌ها دارای هزینه است و محدودیت‌های مشترکی را برای آنها به‌وجود می‌آورد که عبارتند از: هزینه محاسباتی و نیازمندی به حافظه زیاد، بی‌میلی برای استفاده در چند زبان، کمبود پایگاه دانش با کیفیت بالا، برآورد ضعیف از ارتباط معنایی و عدم توانایی برای استفاده در حوزه خاص (محتوا محور) (محمد، ۲۰۰۸). ضعف‌ها و محدودیت‌های سنج‌های تشابه لغوی و معنایی باعث می‌شود که این سنج‌ها میزان

⁴ Imperfection

⁵ Probabilty Theory

⁶ Fuzzy Sets Theory

⁷ Possibilty Theory

⁸ Rough Set Theory

¹ Bloom filters

² Text or Lexical Similarity

³ Semantic Similarity

(جدول - ۱): رویکردها و روش‌های شباهت‌سنجی (محمد، ۲۰۰۸)

نوع تشابه	رویکرد	روش	دقت اندازه‌گیری	حافظه مورد نیاز	زمان پردازش	پیچیدگی محاسبات	پایداری در برابر تغییرات ویرایشی
لغوی	هندسی (بارون و همکاران، ۲۰۰۹)	هم‌پوشانی واژه‌های، سنج‌های فراوانی واژه-وارون فراوانی سند، فاصله ویرایشی، شباهت کسینوسی	متوسط	زیاد	زیاد	کم	کم
	ویژگی محور	مدلهای اثرانگشتی، اثرانگشت غربالی (کونت کنت و همکاران، ۲۰۱۰)، الگوریتم‌های مبتنی بر ترتیب واژه (ژانگ و همکاران، ۲۰۱۱)	خوب	زیاد	زیاد	زیاد	متوسط
	مبتنی بر احتمال (بارون و همکاران، ۲۰۰۹)	ترجمه ماشینی، Okapi BM25، فاصله کالیک-لیپلر	متوسط	زیاد	زیاد	زیاد	خوب
	شبکه استنتاجی (آدی، ۲۰۱۲)	-	بسته به داده آموزشی	زیاد	زیاد	متوسط	متوسط
	یادگیرنده فاصله (بیلنکو و مونی، ۲۰۰۵)	فاصله ویرایشی رشته یادگیری پذیر با شکاف وابستگی، فضای برداری یادگیری پذیر	بسته به داده آموزشی	زیاد	متوسط	متوسط	متوسط
معنایی	مبتنی بر پایگاه دانش (می‌هالیسا و همکاران، ۲۰۰۶)	وو و پالمر، رزینک، هیرست و اس تی آنگ، لی لی کوک و چادرو، لسک	بسته به پایگاه دانش	زیاد	زیاد	زیاد	زیاد
	مبتنی بر پیکره (محمد، ۲۰۰۸)	اطلاعات متقابل نقطه‌ای، تحلیل معنایی پنهان، تحلیل معنایی پنهان احتمالی، تخصیص دیرکله پنهان، مدل ابرفضای قیاس به زبان، نمایه سازی تصادفی	بسته به پایگاه پیکره	زیاد	زیاد	زیاد	کم
	روش‌های هوشمند	شبکه‌های عصبی (یح و همکاران، ۲۰۱۱)، مجموعه فازی (گوپتا و همکاران، ۲۰۱۴)، نظریه گواه (لی و همکاران، ۲۰۰۶)	خوب	زیاد	زیاد	زیاد	زیاد

(جدول - ۲): دسته‌بندی سامانه‌های تشخیص دستبرد ادبی

نام	رویکرد شباهت‌سنجی	نوع تشابه	روش مقایسه
منبر (عثمان و همکاران، ۲۰۱۳)	هندسی	لغوی	واژه به واژه
کپس (عثمان و همکاران، ۲۰۱۳)	هندسی	لغوی	اثرانگشت
اسکم (رام و همکاران، ۲۰۱۴)	هندسی	لغوی	واژه به واژه (در سطح جمله)
چک (کایا و همکاران، ۲۰۱۴)	ویژگی محور	لغوی	ویژگی‌های جمله
رونوشت‌برداری از کالاهای دیجیتال (باو و همکاران، ۲۰۰۱)	هندسی	لغوی	مبتنی بر بردار
اسپلت (کولبرگ و همکاران، ۲۰۰۳)	هندسی	لغوی	تعداد تکرار واژه‌ها
نشان‌گذاری شفاف دیجیتال (زینی و همکاران، ۲۰۱۳)	ویژگی محور	لغوی	اضافه کردن اطلاعات به سند
زینی (زینی و همکاران، ۲۰۱۳)	ویژگی محور	لغوی/معنایی	خوشه‌بندی
چاو و رحمان (چاو و رحمان، ۲۰۰۹)	ویژگی محور	لغوی/معنایی	شبکه‌های عصبی
اقبال (اقبال و همکاران، ۲۰۱۲)	هندسی	لغوی	مبتنی بر بردار و گراف
گراوند (گراوند و همکاران، ۲۰۱۴)	هندسی	لغوی	استفاده از فیلترهای بلوم

و باید چارچوب توسعه پیدا کند تا با دانش اضافه شده عناصر $\theta_{N+2}, \theta_{N+1}$ و ... تطبیق یابد (گیارانتانو، ۲۰۰۵).

۴-۲- تابع انتساب باور پایه

مقداردهی در نظریه گوه با «تابع انتساب باور پایه» انجام می شود. این شاخص را می توان با استفاده از نگاشتی مانند m برای بیان باور خود درباره یک گزاره، به شکل عددی در بازه $[0, 1]$ ارائه داد (خطیبی و منتظر، ۲۰۱۰):

$$(2) \quad m: \mathcal{P} \rightarrow [0, 1]$$

این تابع دارای خواص زیر است:

الف- برای هر یک از گزاره ها، باور می تواند مقداری بزرگ تر یا مساوی صفر داشته باشد؛ یعنی:

$$(3) \quad \forall A \in \mathcal{P} \rightarrow m(A) \geq 0$$

ب- هیچ درجه ای از باور برای گزاره تهی در نظر گرفته نمی شود؛ یعنی:

$$(4) \quad m(\emptyset) = 0$$

ج- مجموع کل باور برابر با واحد است؛ یعنی:

$$(5) \quad \sum_{A \in \mathcal{P}} m(A) = 1$$

۴-۳- تابع باور و تابع مقبولیت

«تابع باور» نشان دهنده درجه باور کامل به گزاره ای است که نسبت به آن اطمینان وجود دارد. در این تابع، تنها گزاره هایی وارد می شوند که به شکل کامل مؤید گزاره مدنظر هستند. این تابع به صورت رابطه (۶) تعریف می شود (بای، ۲۰۰۴):

$$(6) \quad Bel(A) = \sum_{C_i \in \mathcal{P}, C_i \subset A} m(C_i)$$

که $m(C_i)$ بخشی از باور کامل منتسب به گزاره A است.

«تابع مقبولیت» نیز بیان گر درجه مقبولیت یک گزاره است. در واقع درجه مقبولیت یک گزاره، به معنای حد بیشینه وقوع آن است. چون هر گزاره ای که با گزاره مدنظر اشتراک دارد، حداقل بخشی از آن گزاره را پوشش می دهد و بدین ترتیب درجه باور گزاره دارای اشتراک بر امکان پذیری و مقبولیت گزاره مدنظر دلالت می کند. این تابع به صورت رابطه (۷) تعریف می شود (گوآن و همکاران، ۱۹۹۱):

$$(7) \quad Pls(A) = \sum_{C_i \cap A \neq \emptyset} m(C_i)$$

در مسائل دنیای واقعی بخشی از نقص اطلاعات و داده ناشی از شناخت ناقص فضای حالت است به طوری که عوامل اثرگذار بر مسئله و همچنین فضای حالت آن به راحتی قابل شناسایی نیست. در این موارد، عدم اطمینانی نسبت به فضای حالت مسئله به وجود می آید که با روش های احتمالی قابل بررسی نیست. این عدم اطمینان نوعی عدم قطعیت به مسئله وارد می کند. به عنوان مثال، آیا با فرض دانستن اینکه به احتمال ۱۰٪ در حیاط دانشگاه نفت وجود ندارد، می توان به این اطمینان رسید که به احتمال ۹۰٪ در حیاط دانشگاه نفت وجود دارد؟ آرتور دمپستر^۱ نخستین بار در دهه ۱۹۶۰ به این موضوع اشاره کرد و تلاش کرد تا به جای استفاده از یک عدد احتمالی منفرد برای نمایش عدم قطعیت؛ آن را به صورت بازه ای از احتمالات مدل کند تا بتواند شناخت ناکاملی را از فضای حالت مسئله بیان نماید (دمپستر، ۱۹۶۷). شفر^۲ ضمن اصلاح و گسترش کار دمپستر، کنلی را به نام «نظریه ریاضی گوه»^۳ در سال ۱۹۶۷ منتشر کرد (شفر، ۱۹۷۵). کار گسترده تر شفر «استدلال مبتنی بر گوه»^۴ نام داشت که به تفسیر اطلاعات غیرقطعی، نادقیق و برخی ناصحیح می پرداخت (گریور، ۱۹۸۱).

اجزای اصلی این نظریه گوه شامل: چارچوب تشخیص^۵، توابع انتساب باور پایه^۶، باور^۷ و مقبولیت^۸ است که در ادامه توضیح داده می شود:

۴-۱- چارچوب تشخیص

«چارچوب تشخیص»، مجموعه ای از گزاره های کامل و منحصر به فرد (همانند فضای حالت در نظریه احتمال در بررسی مسئله است) که چارچوبی را برای پوشش همه حالت های مسئله ارائه می کند. مجموعه چارچوب تشخیص را با نماد Θ و با رابطه ۱ نمایش می دهند (گوآن و بل، ۱۹۹۱):

$$(1) \quad \Theta = \{\theta_1, \theta_2, \dots, \theta_N\}$$

که θ_i حالت های ممکن رخداد در فضای حالت را نشان می دهد و منظور از «تشخیص» این است که امکان تشخیص یک پاسخ صحیح از میان همه پاسخ های ممکن وجود دارد. اگر پاسخ در چارچوب نباشد، نشان دهنده آن است که شناخت کاملی از فضای حالت مسئله وجود نداشته

¹ Arthur Dempster

² Shafer

³ A mathematical theory of evidence

⁴ Evidential reasoning

⁵ Frames of discernment

⁶ Basic probability assignment function

⁷ Belief function

⁸ Plausibility function

در این رابطه m_i تابع انتساب گواه پایه گزاره و w_i وزن این انتساب (بر اساس اعتبار اطلاعات موجود) است.

۵- معماری سامانه تشخیص دستبرد ادبی

شکل (۳) معماری سامانه پیشنهادی را نشان می‌دهد. در این معماری مراحل زیر به ترتیب بر روی دو سند مشکوک و مرجع انجام خواهد شد:

۱- قطعه‌بندی و ثبت سند در پایگاه داده

۱-۱- قطعه‌بندی کردن جمله‌های اسناد

براساس فهرستی از علائم معیار همانند

نقطه‌گذاری و نقطه‌ویرگول

۱-۲- ثبت جمله‌ها و ترتیب آنها در پایگاه داده

۱-۳- قطعه‌بندی جمله‌ها براساس معیار

«فاصله» بین واژه‌ها

۱-۴- ثبت واژه‌ها به‌همراه ترتیب آنها برای

ساخت چندکلمه‌ها در پایگاه داده

۱-۵- تولید چندکلمه‌های واژه برای تولید

عبارت چندواژه‌ای

۲- پیش‌پردازش سند

۲-۱- پیش‌پردازش واژه‌ها و عبارت چندواژه‌ای

و حذف ایستواژه‌ها (مانند «ز» و «به»)

۲-۲- ریشه‌یابی واژه‌ها و عبارت چندواژه‌ای

(به‌طورمثال ریشه واژه «خور» «می‌خورند»

است.)

۳- تحلیل سند

۳-۱- ساخت ماتریس تعداد تکرار واژه و

اهمیت واژه (TF-IDF) از روی واژه‌ها و

عبارت چندواژه‌ای ساخته شده

۳-۲- یافتن واژه‌ها و عبارت چندواژه‌ای اصلی

براساس تعداد وقوع

۳-۳- جداسازی واژه‌های تخصصی براساس

هستان‌نگار موضوعی

۳-۴- تشکیل ماتریس جمله‌های تخصصی و عمومی

۴- محاسبه شباهت بین اسناد

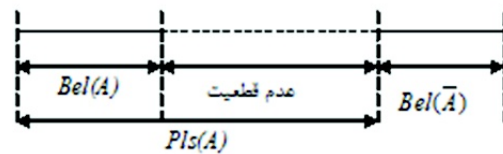
۴-۱- اندازه‌گیری میزان شباهت زوج جمله‌ها بر

اساس ماتریس جمله‌های عمومی دو سند

با استفاده از واژه‌نامه

در صورت فقدان اطلاعات و دانش کافی درباره مسئله،

بسیار معقول است که به جای ارایه یک مقدار معین برای متغیر نامطمئن (همانند آنچه که در نظریه احتمال انجام می‌شود) بازه‌ای برای محدوده عدم قطعیت ارائه شود. با توجه به تعریف تابع باور، این تابع حداقل درجه گواه کامل گزاره را که درباره آن اطمینان وجود دارد، ارائه می‌کند؛ اما تابع مقبولیت، حداکثر درجه باور بالقوه گزاره را بیان می‌کند. بدین ترتیب، مقدار متغیر عدم قطعیتی (درجه باور) بین دو کران پایین و بالا، که به ترتیب تابع باور و تابع مقبولیت است، محصور می‌شود. همان طور که در شکل ۲ ملاحظه می‌شود، درجه باور گزاره فرضی A با بازه $[Bel(A), Pls(A)]$ ارائه می‌شود. با توجه به اینکه مقادیر توابع باور و مقبولیت در بازه $[0, 1]$ هستند، درجه باور گزاره نیز در بازه $[0, 1]$ جای دارد (خطیسی و منتظر، ۲۰۱۰).



(شکل-۲): توابع باور و مقبولیت یک گزاره

۴-۴ ترکیب گواه

اگر چند تابع انتساب گواه پایه برای چارچوب تشخیص مسئله ارائه شود، باید آنها را با هم ترکیب کنیم، تا تابع انتساب جدیدی حاصل شود که حاوی معتبرترین و کامل‌ترین اطلاعات باشد (گوان و همکاران، ۱۹۹۱). به لحاظ ریاضی، دو گواه m_1 و m_2 را که از طریق دو منبع اطلاعاتی مستقل حاصل آمده، می‌توان با استفاده از قاعده ترکیب گواه ترکیب کرد و به تابع انتساب جدیدی دست یافت؛ این ترکیب با رابطه (۸) تعریف می‌شود:

$$m(A) = \frac{\sum_{C_i \cap C_j = A} m_l(C_i) m_r(C_j)}{1 - \sum_{C_i \cap C_j = \emptyset} m_l(C_i) m_r(C_j)}, \quad A \neq \emptyset \quad (8)$$

در این قاعده، C_i و C_j هر یک به گزاره‌هایی از یک منبع

مستقل اشاره می‌کنند. روش میانگین‌گیری^۱ صحت اطلاعات ساختارهای گواه را یکسان فرض می‌کند و از تناقض موجود بین ساختارهای گواه چشم‌پوشی می‌کند. قاعده ترکیب میانگین‌گیری به‌صورت رابطه (۹) است (بای، ۲۰۰۴):

$$m_{l..n} = \frac{1}{n} \sum_{i=1}^n w_i m_i \quad (9)$$

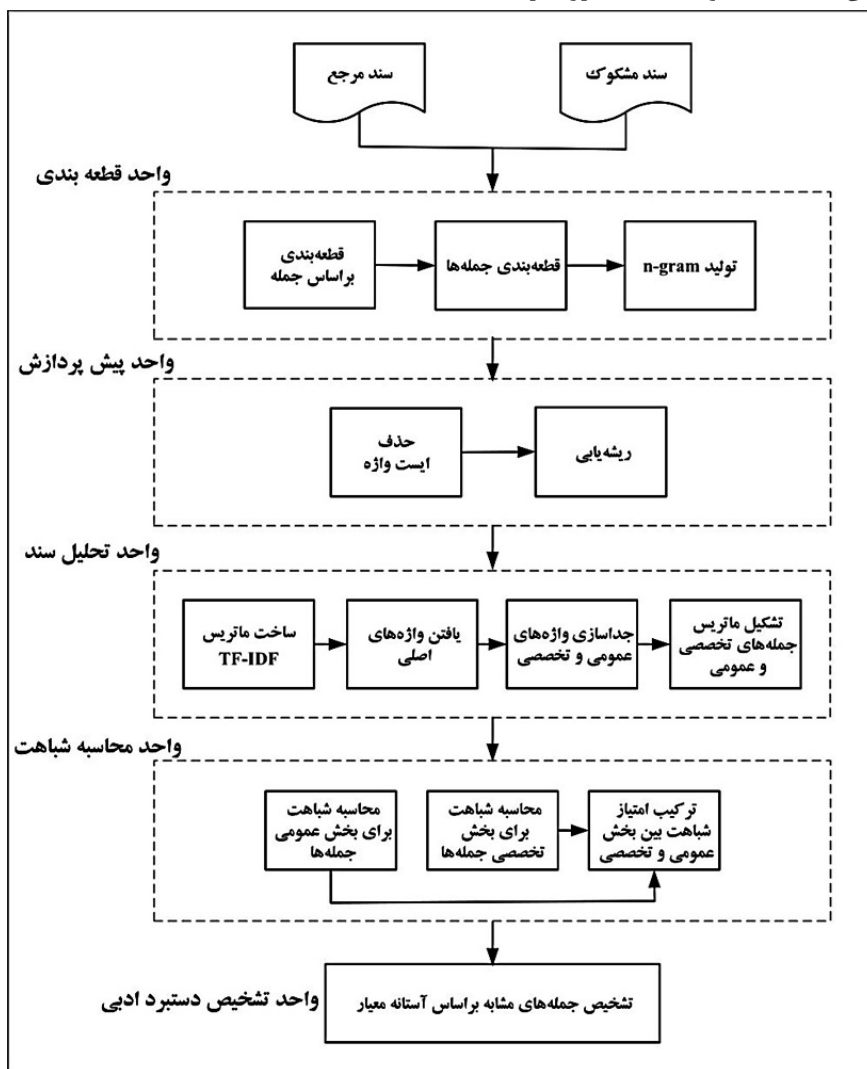
¹ Averaging Method

داخلی، به شدت وابسته به محتوای متن و ساختار طبقه‌بندی واژه‌نامه (مانند ارتفاع، ساختار و زیردرخت) هستند و کارایی و سودمندی آن با توجه به طبقه‌بندی و ساختار مفاهیم از هم متفاوت و گوناگون خواهد بود (پدرسون و همکاران، ۲۰۰۴؛ ووهرینگر و فلیدل، ۲۰۱۱) به همین دلیل این سنج‌ها سنجش نادقیق و ناکاملی از زوج واژه‌ها به دست می‌دهند. از همین رو، در این مقاله، از چارچوب نظریه گواه برای ترکیب سنج‌های مختلف استفاده شده است. در این ایده، هر سنج به‌عنوان یک گواه اطلاعاتی در نظر گرفته شده و سپس از نظریه گواه برای پوشش معایب هر سنج استفاده و با اعمال قانون ترکیب گواه‌ها، امتیاز نهایی برای شباهت‌سنجی دو جمله به دست می‌آید.

۲-۴- اندازه‌گیری میزان شباهت زوج جمله‌ها براساس ماتریس جمله‌های تخصصی دو سند با استفاده از هستان‌نگار موضوعی
 ۳-۴- ترکیب امتیاز شباهت عمومی و تخصصی برای زوج جمله‌ها و محاسبه امتیاز نهایی بین زوج جمله‌ها

۵- تشخیص جمله‌های مشابه

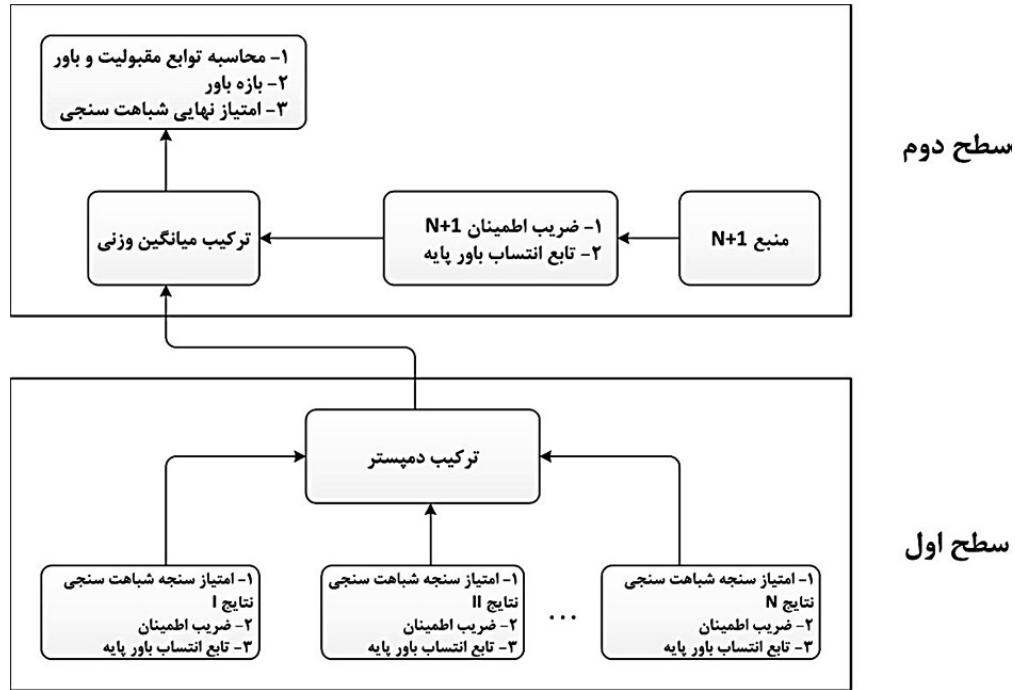
باید توجه داشت که هسته اصلی و قلب این سامانه «واحد محاسبه شباهت» است که براساس آن پس از انجام پردازش‌های نخستین بر محتوای سند، با استفاده از سنج‌های شباهت‌سنجی، میزان تشابه آنها اندازه‌گیری می‌شود. پژوهش‌ها نشان می‌دهد که هیچ‌یک از سنج‌های شباهت‌سنجی در حالت کلی برتر از بقیه نیست (محمد، ۲۰۰۸). این سنج‌ها با توجه به سازوکار



(شکل - ۳): معماری سامانه تشخیص دستبرد ادبی

استفاده) تخصیص داده می شود. در سطح دوم نتیجه حاصل از سطح نخست با باور پایه حاصل از گواه سوم از طریق تابع ترکیب میانگین گیری (رابطه ۹) با هم ترکیب و سپس این نتیجه برای محاسبه تابع باور و مقبولیت (روابط ۶ و ۷) استفاده می شود. این توابع به عنوان بازه تصمیم گیری برای چارچوب تشخیص در حل مسئله در نظر گرفته می شود و در نهایت تصمیم نهایی بر روی بهترین بازه انتخاب می شود و شباهت بین دو جمله به دست می آید.

شکل (۴) جزییات انجام این فرآیند را نشان می دهد. در این ساختار، n سنجه در سطح نخست به عنوان گواه در نظر گرفته می شوند و سپس امتیاز حاصل از سنجش شباهت بین دو جمله این n سنجه، به عنوان باور پایه برای چارچوب تشخیص نظریه گواه در نظر گرفته می شود. این نتایج از طریق تابع هم جوشی دمپستر (رابطه ۸) با هم ترکیب شده و به عنوان باور پایه به سطح دوم فرستاده می شود. در این مرحله برای به دست آوردن باور کافی نسبت به نتایج سطح قبلی، ضریب اطمینانی که از طریق کارشناسان و خبرگان منتخب تعیین شده به هر یک از گواهان (سنجه های مورد



(شکل - ۴): فرآیند هم جوشی سنجه های مختلف در بخش محاسبه شباهت سنجی

در این سامانه، جمله ها به دو بخش عمومی و تخصصی بر اساس واژه های موجود در پایگاه دانش و هستان نگار^۱ تقسیم و سپس زوج واژه های آنها با هم مقایسه می شوند. برای بخش عمومی از امتیاز ارتباط لغوی بین واژه ها با استفاده از واژه نامه عمومی استفاده شده به صورتی که چنانچه دو واژه در فهرست مترادف های هم باشند امتیاز یک و در غیر این صورت امتیاز صفر در نظر گرفته می شود و برای بخش تخصصی از روش های مبتنی بر پایگاه دانش استفاده می شود. در بخش تخصصی واژه های متن به عنوان مفهوم^۲ فرض می شوند. این مفاهیم به طور معمول در هستان نگارها و شبکه های معنایی^۳، که به صورت

در این سامانه، جمله ها به دو بخش عمومی و تخصصی بر اساس واژه های موجود در پایگاه دانش و هستان نگار^۱ تقسیم و سپس زوج واژه های آنها با هم مقایسه می شوند. برای بخش عمومی از امتیاز ارتباط لغوی بین واژه ها با استفاده از واژه نامه عمومی استفاده شده به صورتی که چنانچه دو واژه در فهرست مترادف های هم باشند امتیاز یک و در غیر این صورت امتیاز صفر در نظر گرفته می شود و برای بخش تخصصی از روش های مبتنی بر پایگاه دانش استفاده می شود. در بخش تخصصی واژه های متن به عنوان مفهوم^۲ فرض می شوند. این مفاهیم به طور معمول در هستان نگارها و شبکه های معنایی^۳، که به صورت

۵-۱- روش وو و پالمر^۴

در این روش تشابه معنایی با توجه به عمق (ارتفاع) دو مفهوم مورد مقایسه در طبقه بندی ساختار شبکه معنایی و همچنین مقدار عمق مفهوم مشترک بین این دو مفهوم در کمترین فاصله به دست می آید:

$$Sim(C_1, C_2) = \frac{\bullet \cdot depth(LCS(C_1, C_2))}{depth(C_1) + depth(C_2)} \quad (10)$$

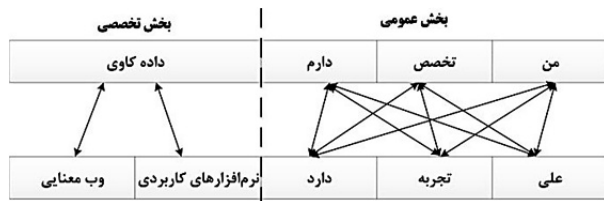
⁴ Wu and Palmer

¹ Ontology

² Concept

³ Semantic networks

کسب نتیجه بهتر، هر یک از سنجه‌های مورد استفاده را به‌عنوان یک منبع (گواه) در نظر می‌گیریم و نتایج آنها را با استفاده از نظریه گواه ترکیب می‌کنیم. همچنین در این مرحله با استفاده از هستان‌نگار، واژه‌های تخصصی را از واژه‌های عمومی جدا و هر جمله را به دو شبه‌جمله تقسیم می‌کنیم. برای نمونه دو جمله «من در داده کاوی تخصص دارم» و «علی در نرم‌افزارهای کاربردی و وب‌معنایی تجربه دارد» را در نظر بگیرید؛ این جمله‌ها با استفاده از هستان‌نگار و واژه‌نامه، مطابق با شکل (۵) به دو بخش عمومی و تخصصی تقسیم و زوج‌واژه‌های آنها با هم مقایسه می‌شوند. در این مقاله از هستان‌نگار تخصصی فناوری اطلاعات (حورعلی، ۱۳۹۰) استفاده شده است. یادآور می‌شود چنانچه ارتباط بین زوج‌واژه‌های در هستان‌نگار تعریف نشده باشد؛ ابتدا از واژه‌نامه عمومی برای مقایسه آنها و سپس پیشنهادی به کاربر سامانه برای گسترش هستان‌نگار داده خواهد شد.



(شکل - ۵): نحوه مقایسه زوج‌واژه‌ها در زوج جمله‌ها

در این مقاله از روابط ۱۰ تا ۱۲ برای محاسبه شباهت بین بخش تخصصی دو جمله استفاده شده است. در جدول ۳ مقدار شباهت زوج‌جمله بالا از طریق روابط بالا آورده شده است:

(جدول - ۳): امتیاز شباهت بین دو جمله براساس سنجه‌های مختلف

نام سنجه	امتیاز شباهت بین دو جمله
وو و پالمر (رابطه ۱۰)	۰/۷۱
رزنیک (رابطه ۱۱)	۰/۹۰
رزنیک مبتنی بر تعداد گره (رابطه ۱۲)	۰/۸۴

با توجه به جدول (۳)، ملاحظه می‌شود، سنجه رزنیکی هم‌خوانی بهتری با نظرهای واقعی دارد و به همین دلیل برای بخش تخصصی از آن استفاده شده است. پس از

LCS^۱ مفهوم مشترک بین دو مفهوم C_1, C_2 در ساختار درختی مفاهیم در کمترین فاصله درخت است و تابع $depth$ عمق مفهوم (ارتفاع در درخت) این مفاهیم را به‌دست می‌آورد (وو و همکاران، ۱۹۹۴).

۵-۲- روش رزنیکی^۲

این سنجه مبتنی بر محتوای اطلاعاتی است که مفهوم مشترک دو واژه در کمترین فاصله را در بردارد و به‌صورت رابطه (۱۱) تعریف می‌شود:

$$rel(w_1, w_2) = \max[rel(C_1, C_2)] \quad (11)$$

$$C_1 \in s(w_1), C_2 \in s(w_2)$$

که در آن $s(w_i)$ مجموعه مفاهیمی متناظر با معانی مختلف واژه w_i است. به عبارتی دیگر، C_1, C_2 مرتبط‌ترین مفاهیم در سلسله‌مراتب شبکه واژگان هستند که جزو مجموعه معانی w_1, w_2 هستند. این سنجه را می‌توان به‌صورت دیگری بر مبنای محتوای اطلاعاتی تعداد فرزندان هر مفهوم در هستان‌نگار نیز تعریف کرد، که دید عینی و واقعی‌تری به دست می‌دهد. نحوه محاسبه این سنجه عبارت از:

$$IC_{onto}(C) = 1 - \frac{\log(num_desc(C)+1)}{\log(max_{onto})} \quad (12)$$

که $num_desc(c)$ تعداد فرزندان هر مفهوم و max_{onto} تعداد مفهوم موجود در هستان‌نگار بوده که برای هنجارسازی مقدار تشابه به کار گرفته می‌شود (رزنیکی، ۱۹۹۵).

۶- پیاده‌سازی عملیاتی سامانه

سامانه پیشنهادی در محیط نرم‌افزاری Visual Studio .Net و با پایگاه داده Microsoft SQL پیاده‌سازی شده و معماری طراحی سامانه به‌صورت شیء‌گرا^۳ انجام گرفته که توانایی سنجش هر دو جمله را براساس سنجه‌های متفاوت و همچنین پایگاه دانش‌های مختلف داشته باشد. برای استفاده از نظریه گواه ابتدا چارچوب تشخیص را به‌صورت فضای حالت شباهت بین دو جمله تعریف می‌کنیم و از آنجا که سنجه‌های مبتنی بر پایگاه دانش از منظرهای مختلف و با کیفیت‌های متفاوتی شباهت بین دو جمله را می‌سنجند، باعث نوعی عدم شناخت کامل در مسئله می‌شود؛ این موضوع استفاده از نظریه گواه را توجیه‌پذیر می‌کند. برای

¹ Least Common Subsumer

² Resnik

³ Object Oriented

فناوری اطلاعات، موجب پدید آمدن نوع جدیدی از یادگیری به نام یادگیری الکترونیکی شده است» جمله مشابه « یادگیری الکترونیکی نوع جدیدی از یادگیری است که با استفاده از ابزارهای فناوری اطلاعات ارائه می شود» ساخته شده است. در نهایت برای ارزیابی دقیق سامانه‌ای، میزان شباهت زوج جمله‌ها توسط تعدادی خبره متفاوت با خبرگان تولیدکننده مشخص شده و امتیاز نهایی شباهت بین دو جمله براساس میانگین امتیاز خبرگان مدنظر قرار گرفته است.

از آنجا که هدف در این مقاله بررسی شباهت در متون تخصصی فارسی است و در حال حاضر هیچ‌گونه پیکره یا مجموعه داده‌ای با این ویژگی در زبان فارسی و همین‌طور زبان‌های دیگر وجود ندارد، از مجموعه پایگاه مقاله‌های کنفرانس یادگیری الکترونیکی برای تولید پیکره استفاده شده است. این پیکره شامل ۸۱۰ زوج جمله است که جمله نخست زوج جمله‌ها از مقاله و جمله دوم توسط خبرگان زبان به صورت تصادفی در میزان شباهت‌های مختلف تولید شده است. به عنوان نمونه برای جمله «گسترده‌گی تحولات در

(جدول - ۵): ترکیب نتایج با استفاده از نظریه گواه

باور پایه				ب.ب.ب.ب.	گواهان
ناباوری	نامشابه	نیمه مشابه	مشابه		
۰/۰۲۹	۰	۰	۰/۹۷۱	۰/۵	نتیجه سطح اول (روابط ۱۰ و ۱۱)
۰/۱۶	۰	۰	۰/۸۴	۰/۵	نتیجه سطح دوم (رابطه ۱۲)
۰/۰۹۴۵	۰	۰	۰/۹۰۵۵	-	ترکیب باورهای پایه دو گواه از طریق قاعده میانگین‌گیری
دو جمله مشابه هستند					نتیجه نهایی

دستبرد ادبی به صورت برون‌خط انجام می‌شود، در سامانه‌های تشخیص دستبرد ادبی دقت تشخیص بیش از سرعت مدنظر است.

(جدول - ۶): نتایج کارایی روش پیشنهادی و مقایسه با سه روش دیگر

روش اجرا (دقیقه)	بازخوانی	دقت	معیار F	روش
۶/۳	۰/۱۸۷	۰/۸۴	۰/۱۸۶	روش گواه
۵/۴	۰/۱۷۲	۰/۶۹	۰/۱۷۱	روش ارائه شده در (موهلر و همکاران، ۲۰۰۹)
۷/۱	۰/۶۴	۰/۶۹	۰/۶۶	روش ارائه شده در (گوپتا و همکاران، ۲۰۱۴)
۵/۱	۰/۶۲	۰/۵۹	۰/۶۰	روش مستقیم (فاصله ویرایشی)

۷- نتیجه‌گیری و پژوهش‌های آتی

در این مقاله از نظریه گواه برای هم‌جوشی اطلاعات به منظور شباهت‌سنجی بین دو سند و کشف دستبرد ادبی، با توجه به کیفی و ناکامل بودن عوامل اثرگذار سنجش شباهت بین دو متن استفاده شده است. سامانه طراحی شده با استفاده از

تاکنون هیچ‌گونه پیکره تخصصی برای تشخیص شباهت‌سنجی در متون تخصصی ارائه نشده است و بیش‌تر روش‌های ارائه شده در ادبیات بر روی پیکره‌های انگلیسی انجام می‌پذیرد و نمی‌توان آنها را با روش پیشنهادی مقایسه کرد. از این رو، در این مقاله نتایج پیاده‌سازی روش‌های ارائه شده در مقاله‌های (گوپتا و همکاران، ۲۰۱۴؛ موهلر و همکاران، ۲۰۰۹)، فاصله ویرایشی به عنوان روش مستقیم و روش پیشنهادی بر روی پیکره فارسی تخصصی تولید شده، پیاده‌سازی شده و نتایج آنها بر مبنای معیارهای F، دقت، بازخوانی و زمان با هم مقایسه شده‌اند. نتایج این مقایسه در جدول ۶ ارائه شده است. نتایج بیان‌گر آن است که روش پیشنهادی به میزان قابل قبولی، بهتر از سه روش دیگر بوده و توانایی شناسایی متون مشابه را دارد.

با توجه به نتایج جدول (۶) ملاحظه می‌شود استفاده از معماری پیشنهادی و استفاده از هم‌جوشی سنج‌ها با کمک نظریه گواه، باعث بهبود دقت در عملکرد سامانه شده و همچنین بازخوانی اسناد مرتبط و مشابه را بهتر انجام می‌دهد. از طرفی در بُعد سرعت از آنجا که در روش پیشنهادی از محاسبات بیشتری استفاده می‌شود نسبت به دیگر روش‌ها تفاوت خاصی دیده نمی‌شود. اجرای این الگوریتم‌ها بر روی رایانه‌ای با CPU P4 و RAM 4g انجام شده است. لازم به ذکر است با توجه به اینکه تشخیص

Conference on Knowledge Management in Organizations, pp. 285-297. Springer Netherlands.

Bae, h.-r., Grandhi, r. V. & Canfield, r. A. (2004). Epistemic Uncertainty quantification techniques including evidence theory for large-scale structures. *Computers and Structures*, 23, 125-138.

Bao J, Shen J, Liu X. (2001). On illegal copying and distributing detection mechanism for digital goods. *J Comput Res Develop*; 38(1):121-5.

Barron-Cedeno, A. Eiselt, A. Rosso, P. Monolingual Text Similarity Measures: A Comparison of Models over Wikipedia Articles Revisions, Proc. 7th Int. Conf. on Natural Language Processing, ICON-, Hyderabad, India, pp. 29-38, 2009.

Basile, C., Benedetto, D., Caglioti, E., Cristadoro, G., & Esposti, M. D. (2009). A plagiarism detection procedure in three steps: Selection, matches and "squares". In Proceedings of the 3rd PAN Workshop. *Uncovering Plagiarism, Authorship and Social Software Misuse*.

Bilenko Mikhail and Mooney Raymond J. Adaptive Duplicate Detection Using Learnable String Similarity Measures. Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD-2003), Washington DC, pp.39-48, August, 2003.

Bretag, T., & Mahmud, S. (2009). Self-plagiarism or appropriate textual re-use? *Journal of Academic Ethics*, 7, 193-205.

Chow, K.K. Salim, N. (2010). Web based cross language plagiarism detection, in: Second International Conference on Computational Intelligence, Modelling and Simulation, pp. 199-204.

Collberg C, Kobourov S, Louie J, Slattery T. (2003). Splat: a system for self-plagiarism detection. In: Proceedings of IADIS international conference WWW/INTERNET; p. 508-14.

Dempster, A. (1967). Upper and Lower Probabilities Induced by Multivalued Mapping, *Annals of Math. Stat.*, 38, 325-329.

Devi, S. L., Rao, P. R. K., Ram, V. S., & Akilandeswari, A. (2010). External plagiarism detection - Lab report for PAN at CLEF 2010. In Notebook Papers of CLEF LABs and Workshops.

Dreher, H. (2007). Automatic conceptual analysis for plagiarism detection. *Information and Beyond: The Journal of Issues in Informing Science and Information Technology*, 4, 601-614.

Ekbal A, Saha S, Choudhary G. (2012). Plagiarism detection in text using vector space model. In: 12th International conference on hybrid intelligent systems (HIS); p. 366-71.

Frantzi, K., Ananiadou, S., Mima, H. (2000). Automatic recognition of multi-word terms. *International Journal of Digital Libraries*, 3, 2, 115-130.

Geravand S, Ahmadi M. (2014). An efficient and scalable plagiarism checking system using Bloom filters. *Comput Electr Eng*.

گواهای موجود، استنتاج را در دو سطح انجام می‌دهد. در مرحله ابتدایی جمله‌های موجود در سند به دو بخش عمومی و تخصصی تقسیم و سپس با استفاده از سنج‌های متفاوت و همچنین استفاده از هستان‌نگار تخصصی امتیاز تشابه برای هر بخش را محاسبه کرده و در نهایت امتیاز تجمعی از دو بخش برای دو جمله به دست آورده می‌شود در سطح نخست نتایج سنج‌های شباهت‌سنجی به‌عنوان گوا (با باور پایه مشخص) با قاعده ترکیب دمپرستر-شفر با هم ترکیب و به‌عنوان گواهی جدید به سطح دوم منتقل می‌شوند. در سطح دوم، نتیجه سطح نخست و گواهی جدید از طریق قاعده میلیگین‌گیری هم‌جوش شده و توابع باور و مقبولیت نهایی محاسبه و شباهت بین دو جمله ارزیابی می‌شود.

این سامانه برای شباهت‌سنجی در پیکره تولیدشده توسط پژوهش‌گران که شامل ۸۱۰ زوج جمله است به کار گرفته شده است. با توجه به نتایج حاصل، ملاحظه شد سامانه در بیش از ۸۶٪ موارد، تحلیل درستی از شباهت دو جمله ارائه می‌دهد و می‌توان از این سامانه برای شناسایی مصادیق دستبرد ادبی استفاده کرد.

در حال حاضر سامانه پیشنهادی با توجه به تعداد کم اسناد با سرعت خوبی پاسخ‌گوست و انتظار این است که با افزایش حجم اسناد، سامانه به سمت کندی رود که می‌توان از روش‌هایی همچون اثرانگشتی و خوشه‌بندی برای کاهش حجم سندهای مورد مقایسه استفاده کرد. از طرفی پایگاه دانش به کار گرفته شده باید غنی شود تا توانایی شناسایی ارتباط واژه‌های جدید را داشته باشد. ترکیب سنج‌های مختلف با استفاده از روش‌هایی مانند نظریه فازی می‌تواند به بهبود اندازه‌گیری شباهت معنایی و شناسایی روابط معنایی بین واژه‌ها منجر شود. از طرفی می‌توان از روش‌های شبکه‌های عصبی و ترکیب آنها با الگوریتم‌های ژنتیک برای تشخیص عبارت‌های چندواژه‌ای نیز استفاده کرد تا میزان دقت مقایسه واژه‌ها را بهبود داد.

۸- مراجع

حورعلی، مریم. یادگیری هوشمند هستان‌نگار برای بسط پرسرمان در جستجوی معنایی، رساله دکتری، دانشگاه تربیت مدرس، ۱۳۹۰، صص ۳۰-۴۵.

Addis, Andrea. Study and Development of Novel Techniques for Hierarchical Text Categorization. Department of Electrical and Electronic Engineering. University of Cagliari, 2010.

Alfred, Rayner, Leow Ching Leong, Chin Kim On, and Patricia Anthony. (2014) "A Literature Review and Discussion of Malay Rule-Based Affix Elimination Algorithms." In The 8th International

- M., Mohler, and Rada Mihalcea. "Text-to-text semantic similarity for automatic short answer grading." Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics. Association for Computational Linguistics, 2009.
- Metzler, D. , Bernstein, Y. , Croft, W. B. , Moffat, A. , and Zobel, J. (2005). Similarity measures for tracking information flow. In Proceedings of CIKM '05, pp. 517-524.
- Metzler, D. Dumais, S. Meek, C. (2007). Similarity Measures for Short Segments of Text , In Proc. of ECIR-07 Springer, Vol. 4425 , pp. 16-27 ,2007.
- Meuschke, N., Gipp, B., (2013). "State-of-the-art in detecting academic plagiarism", International Journal for Educational Integrity Vol. 9 No. 1, pp. 50–71.
- Mihalcea, R. Corley, C. Strapparava, C. (2006). Corpus-based and Knowledge-based Measures of Text Semantic Similarity, In American Association for Artificial Intelligence , pp. 775-780.
- Mohammad, Saif. (2008). Measuring Semantic Distance Using Distributional Profiles Of Concepts, University of Toronto.
- Muhr, M., Zechner, M., Kern, R., & Granitzer, M. (2009). External and intrinsic plagiarism detection using vector space models. In Proceedings of the 3rd PAN Workshop. Uncovering Plagiarism, Authorship and Social Software Misuse, pp.47–55.
- Osman, Ahmed Hamza, Naomie Salim, and Ammar Ahmed E. Elhadi. (2013). "A tree-based conceptual matching for plagiarism detection." In Computing, Electrical and Electronics Engineering (ICCEEE), 2013 International Conference on, pp. 571-579. IEEE.
- Osman, Ahmed Hamza. Salim, Naomie. Binwahlan, Mohammed Salem. Alteeb, Rihab. Abuobieda, Albaraa. (2012). An improved plagiarism detection scheme based on semantic role labeling, Applied Soft Computing 12 , 1493–1502.
- Rahimtoroghi, E.; Faili, H.; Shakery, A, (2010). "A structural rule-based stemmer for Persian," Telecommunications (IST), 5th International Symposium on , vol., no., pp.574,578, 4-6.
- Ram, R. Vijay Sundar, Efstathios Stamatatos, and Sobha Lalitha Devi. (2014). "Identification of Plagiarism Using Syntactic and Semantic Filters." In Computational Linguistics and Intelligent Text Processing, pp. 495-506. Springer Berlin Heidelberg.
- Resnik, P. (1995), Using information content to evaluate semantic similarity. In Proceedings of the 14th International Joint Conference on Artificial Intelligence.
- Shafer, G. (1976). A Mathematical Theory of Evidence, New Jersey, Princeton University Press.
- Sheykh Esmaili, k., Neshati, m., Abolhassani, a. (2006) improving ir performance using intelligent query expansion. 11 international csi computer conferences (csicc'2006).
- Gipp, B., & Beel, J. (2010). Citation based plagiarism detection: A new approach to identify plagiarized work language independently. In Proceedings of the 21st ACM Conference on Hypertext and Hypermedia (pp. 273–274).
- Givi, H.A, Anvari, H. (2006). "Persian Language", 27th ed., Fatemi, Tehran.
- Guan, J. W., & Bell, D. A. (1991). Evidence Theory and its applications. Amsterdam: Elsevier Science Publisher B.V.
- Gupta, Rohit, et al. "UoW: NLP techniques developed at the University of Wolverhampton for Semantic Similarity and Textual Entailment." SemEval 2014 (2014): 785.
- Hariharan, S., Kamal, S., Faisal, A. V. M., Azharudheen, S. M., & Raman, B. (2010). Detecting plagiarism in text documents. In Proceedings of the International Conference on Recent Trends in Business Administration and Information Processing: Vol. 70. Communications in Computer and Information Science (pp. 497–500). Trivandrum, Kerala, India: Springer.
- Heather, J. (2010). Turnitoff: Identifying and fixing a hole in current plagiarism detection software. Assessment & Evaluation in Higher Education, 35(6), 647–660.
- Jiang, J. , and Conrath, D. (1997), Semantic similarity based on corpus statistics and lexical taxonomy. In Proceedings of the International Conference on Research in Computational Linguistics, pp 19–33.
- Joy A.M.L.M., (1999). Plagiarism in programming assignments, IEEE Transactions on Education 42 (1) 129–133.
- Khaleghi, Bahador, Alaa Khamis, Fakhreddine O. Karray, and Saiedeh N. Razavi. "Multisensor data fusion: A review of the state-of-the-art." Information Fusion 14, no. 1 (2013): 28-44.
- Khatibi, V. Montazer, G.A. (2010). A fuzzy-evidential hybrid inference engine for coronary heart disease risk assessment ,Expert Systems with Applications, 37(12): p. 8536-8542.
- Kok Kent, C. , Salim, N. , Features Based Text Similarity Detection , Journal Of Computing, Volume 2, Issue 1, pp. 53-57 ,2010.
- Leacock, C. , and Chodorow, M. (1998). Combining local context and WordNet sense similarity for word sense identification. In WordNet, An Electronic Lexical Database. The MIT Press.
- Lesk, M. (1986). Automatic sense disambiguation using machine readable dictionaries: How to tell a pine cone from an ice cream cone. In Proceedings of the SIGDOC Conference, pp. 24–26.
- Li, Yuhua. McLean, David. Bandar, Zuhair A. O'Shea, James D. and Keeley Crockett . Sentence Similarity Based on Semantic Nets and Corpus Statistics. IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 18, NO. 8, AUGUST 2006.



حمید آهانگر بهان دوره کارشناسی خود را در سال ۱۳۸۱ در رشته مهندسی کامپیوتر در دانشگاه شاهد تهران گذرانده و مدرک کارشناسی ارشد خود را در سال ۱۳۸۶ در رشته مهندسی صنایع از دانشگاه شریف اخذ کرده است. ایشان در حال حاضر دانشجوی دوره دکتری مهندسی صنایع در دانشگاه تربیت مدرس است. زمینه‌های علمی مورد علاقه وی تحلیل داده و سامانه‌های اطلاعاتی با روش‌های نرم رایانش است. نشانی رایانامه ایشان عبارت است از :

Ahangarbahan@gmail.com



غلامعلی منتظر در سال ۱۳۴۸ در کازرون (فارس) به دنیا آمد. او در سال ۱۳۷۰ مدرک کارشناسی خود را در رشته مهندسی برق از دانشگاه صنعتی خواجه نصرالدین طوسی و سپس در سال ۱۳۷۳ و

۱۳۷۷ به ترتیب مدارک کارشناسی ارشد و دکتری خود را در همین رشته از دانشگاه تربیت مدرس اخذ کرد. وی پس از اتمام تحصیلات به عضویت هیأت علمی دانشگاه تربیت مدرس در آمد و در حال حاضر دانشیار مهندسی فناوری اطلاعات در این دانشگاه است. حوزه‌های تخصصی وی شامل نرم رایانش (نظریه مجموعه‌های فازی، شبکه‌های عصبی مصنوعی، نظریه مجموعه‌های نادقیق) و کاربرد آن در سامانه‌های اطلاعاتی (همچون سامانه یادگیری الکترونیکی و سامانه هوشمند حمل و نقل) است. وی تاکنون بیش از ۸۰ مقاله در نشریات معتبر علمی و بیش از ۱۹۰ مقاله در کنفرانس‌های معتبر ملی و بین‌المللی منتشر کرده است. وی علاوه بر این موفق به دریافت جوایز معتبر علمی از جمله «برگزیده جشنواره بین‌المللی خوارزمی»، «برنده کتاب سال دانشگاهی ایران»، «پژوهشگر برگزیده آیسسکو» و «متخصص برجسته فناوری اطلاعات ایران» شده است.

نشانی رایانامه ایشان عبارت است از

montazer@modares.ac.ir

Shivakumar,N. Garcia-Molina H., (1995).SCAM: a copy detection mechanism for digital documents, in: 2nd International Conference in Theory and Practice of Digital Libraries (DL 1995), Austin, TX, June 11–13.

Si, A., Leong, Hong, V., & Lau, R. W. H. (1997).CH-ECK: A document plagiarism detection system. In Proceedings of the ACM Symposium on Applied Computing, pp. 70–77.

Stein, B., Koppel, M., & Stamatatos, E. (Eds.). (2007) Plagiarism Analysis Authorship Identification, and Near Duplicate Detection: Vol. 276. CEUR Workshop Proceedings. CEUR-WS.org. in Proceedings of the SIGIR International Workshop, held in conjunction with the 30th Annual International ACM SIGIR Conference, Amsterdam, Netherlands. 2007.

Stein, Rosso, Stamatatos, Koppel, Agirre (Eds.): (2009).PAN'09, pp. 1-9.

Weber-Wulff, D. (2011).Copy, Shake, and Paste – A Blog about Plagiarism written by a Professor for Media and Computing at the HTW. Online Source. Retrieved October 28, 2012, from: <http://-copy-shake-paste.blogspot.com>.

Weber-Wulff, D. Test cases for plagiarism detection software. (2010).In Proceedings of the 4th International Plagiarism Conference, Newcastle upon Tyne, UK.

Wu, Z., and Palmer, M., (1994), Verb semantics and lexical selection. In Proceedings of the Annual Meeting of the Association for Computational Linguistics.

Yerra., Ng.(2005).A SentenceBased Copy Detection Approach for Web Documents. Fuzzy Systems and Knowledge Discovery.

Yih, Wen-tau. Toutanova, Kristina. Platt, John C. Meek, Christopher. (2011). Learning Discriminative Projections for Text Similarity Measures. Proceedings of the Fifteenth Conference on Computational Natural Language Learning, pages 247–256,Portland, Oregon, USA, 23–24 June 2011.

Yih, Wen-tau. Toutanova, Kristina. Platt, John C. Meek, Christopher. Learning Discriminative Projections for Text Similarity Measures. Proceedings of the Fifteenth Conference on Computational Natural Language Learning, pages 247–256,Portland, Oregon, USA, 23–24 June 2011.

Zhang ,J. Sun,Y. Wang, H. He,Y. Calculating Statistical Similarity between Sentences, Journal of Convergence Information Technology, Volume 6, Number 2. , pp. 22-34,2011.

Zhao Jun, Jin Qian-Li, XU Bo. (2005).Semantic Computation for Text Retrieval. Chinese Journal Of Computers, Vol. 28, No. 12, pp. 2068-2078. 12. (in Chinese).

Zini M, Fabbri M, Moneglia M, Panunzi A. (2006). Plagiarism detection through multilevel text comparison. In: Second international conference on automated production of cross media content for multi-channel distribution (AXMEDIS); p. 181–5.