

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

دانشگاه یزد

دانشکده مهندسی برق و کامپیوتر

گروه مهندسی کامپیوتر

پایان نامه

برای دریافت درجه کارشناسی ارشد

مهندسی کامپیوتر - هوش مصنوعی و رباتیک

ترجمه هوشمند اسامی فارسی در بازیابی اطلاعات بین زبانی

استاد راهنما : دکتر علی محمد زارع بیدکی

استاد مشاور : دکتر علیرضا یاری

پژوهش و نگارش : زهره حق‌اللهی

شهریور ۹۱

تعهدنامه اصالت اثر

اینجانب زهره حق‌اللهی تایید می‌کنم که مطالب مندرج در این پایان‌نامه حاصل کار پژوهشی اینجانب است و به دست آورده‌های دیگران که در این نوشته از آنها استفاده شده است مطابق مقررات ارجاع گردیده است. این پایان‌نامه قبلاً برای احراز هیچ مدرک هم سطح یا بالاتر ارائه نشده است.

کلیه حقوق مادی و معنوی این اثر متعلق به دانشکده برق و کامپیوتر دانشگاه یزد می‌باشد.

نام و نام خانوادگی دانشجو: زهره حق‌اللهی

امضای دانشجو:

تقدیم به پدر، مادر و همسر

تشکر و قدردانی

در ابتدا از استاد محترم آقای دکتر زارع بیدکی که استاد راهنمای اینجانب در انجام این پایان‌نامه بوده‌اند و کمک شایانی در مسیر انجام پروژه داشته‌اند، تشکر و قدردانی ویژه داشته باشم و نیز جناب آقای دکتر یاری که در کسوت استاد مشاور پشتیبان اینجانب بوده‌اند، کمال تشکر را دارم. همچنین از اعضای هیأت داوران که متقبل زحمت داوری پایان‌نامه شدند، نهایت امتنان و سپاسگزاری را دارم.

بخشی از حمایت مالی این پایان‌نامه توسط مرکز تحقیقات مخابرات ایران انجام شده است.

چکیده

با گسترش روزافزون استفاده از اینترنت، کاربران و جستجوگران اطلاعات، دیگر تنها به منابع اطلاعاتی که به زبان آنها نوشته شده اکتفا نمی‌کنند. بنابراین سیستم‌های بازیابی اطلاعات بین زبانی، این نیاز را برای کاربران برآورده می‌سازند. به این صورت که می‌توان تنها با دادن پرس‌وجو در یک زبان، اسناد با زبان‌های مختلف را بازیابی کرد. از بین پرس‌وجوهای کاربران، اسامی افراد از اهمیت خاصی برخوردار هستند. زیرا تعداد زیادی از پرس‌وجوها به اسامی افراد اختصاص می‌یابد. همچنین کاربران زیادی صفحات شخصی خود را به زبان‌های دیگر طراحی می‌کنند. بنابراین در صورت تشخیص این اسامی در پرس‌وجوها و تبدیل آنها به زبان‌های دیگر (که نویسه‌گردانی نام دارد)، می‌توان به صفحات شخصی آنها مستقل از زبانی که هستند، دست یافت.

در این پایان‌نامه، مروری کوتاه بر الگوریتم‌های تشخیص نام، سیستم‌های بازیابی اطلاعات بین زبانی و نیز به طور خاص نویسه‌گردانی کلمات انجام شده است. همچنین چالش‌های مطرح شده در این زمینه و روش عملکرد و نقاط قوت و ضعف آنها مورد بررسی قرار گرفته‌اند. در ادامه روشی برای تشخیص اسامی افراد به زبان فارسی با استفاده از موتور جستجو و همچنین روشی نوین به منظور نویسه‌گردانی این اسامی با استفاده از محتوای صفحات وب فارسی پیشنهاد شده است. مطابق آزمایشات انجام شده با استفاده از روش‌های پیشنهادی، می‌توان به دقت بالایی در زمینه تشخیص و ترجمه اسامی خاص و یا به عبارتی بازیابی اطلاعات بین زبانی اسامی افراد دست پیدا کرد.

واژه‌های کلیدی : موتور جستجو، بازیابی اطلاعات بین زبانی، نویسه‌گردانی، اسامی افراد

فهرست مطالب

عنوان	شماره صفحه
فصل ۱. مقدمه	۱
۱-۱ مقدمه	۲
۱-۲ موتور جستجو و بازیابی اطلاعات	۲
۱-۲-۱ بازیابی اطلاعات	۳
۳-۱ بازیابی اطلاعات بین زبانی	۴
۱-۳-۱ چالش‌های بازیابی اطلاعات بین زبانی	۷
۴-۱ تشخیص نام در پرس‌وجو	۸
۵-۱ نویسه گردانی آوایی	۹
۱-۵-۱ چالش‌های نویسه گردانی آوایی	۱۰
۶-۱ معیارهای ارزیابی	۱۱
۱-۶-۱ معیارهای ارزیابی سیستم تشخیص نام در پرس‌وجو	۱۱
۲-۶-۱ معیارهای ارزیابی سیستم نویسه گردانی	۱۲
۷-۱ اهداف پایان‌نامه	۱۴
۸-۱ ساختار پایان‌نامه	۱۵
فصل ۲. مروری بر کارهای گذشته	۱۶
۱-۲ مقدمه	۱۷
۲-۲ کارهای انجام شده در زمینه تشخیص نام	۱۷
۱-۲-۲ یادگیری با ناظر	۱۸

۱۸.....	یادگیری نیمه ناظر	۲-۲-۲
۲۲.....	یادگیری بدون ناظر	۳-۲-۲
۲۳.....	تحقیقاتی دیگر در تشخیص نام.....	۴-۲-۲
۲۴.....	کارهای انجام شده در زمینه بازیابی اطلاعات بین زبانی.....	۳-۲
۲۶.....	ترجمه پرس و جو.....	۱-۳-۲
۲۷.....	ترجمه اسناد.....	۲-۳-۲
۲۹.....	تکنیک بین زبانی.....	۳-۳-۲
۳۱.....	انطباق واژه‌های کلمات هم‌ریشه.....	۴-۳-۲
	فصل ۳. الگوریتم‌های پیشنهادی: تشخیص اسامی فارسی و نویسه‌گردانی آنها در پرس و جوها	
۴۵.....		
۴۶.....	مقدمه	۱-۳
۴۷.....	الگوریتم تشخیص اسامی فارسی در پرس و جو.....	۲-۳
۴۷.....	طرح مسأله	۱-۲-۳
۴۷.....	الگوریتم پیشنهادی	۲-۲-۳
۵۲.....	ارزیابی الگوریتم.....	۳-۲-۳
۵۴.....	الگوریتم نویسه‌گردانی اسامی فارسی.....	۳-۳
۵۴.....	طرح مسأله	۱-۳-۳
۵۵.....	الگوریتم پیشنهادی	۲-۳-۳
۶۶.....	ارزیابی الگوریتم نویسه‌گردانی.....	۳-۳-۳
۷۱.....	نتیجه‌گیری و کارهای آینده.....	فصل ۴.

۷۲.....	نتیجه‌گیری.....	۱-۴
۷۳.....	کارهای آینده.....	۲-۴
۷۵.....	پیوست آ.....	
۷۶.....	پیوست ب.....	
۷۷.....	پیوست ج.....	
۷۹.....	پیوست د.....	
۸۰.....	واژه‌نامه فارسی به انگلیسی.....	
۸۵.....	واژه‌نامه انگلیسی به فارسی.....	
۹۰.....	منابع.....	

فهرست شکل‌ها

عنوان	شماره صفحه
شکل (۱-۱) نمای کلی سیستم بازیابی	۴
شکل (۲-۱) بازیابی اطلاعات بین زبانی	۵
شکل (۳-۱) چندزبانه بودن محیط وب	۶
شکل (۱-۲) روش‌های انطباق در بازیابی اطلاعات بین زبانی	۲۶
شکل (۲-۲) ترجمه پرس‌وجو در بازیابی اطلاعات بین زبانی	۲۷
شکل (۳-۲) ترجمه اسناد در بازیابی اطلاعات بین زبانی	۲۸
شکل (۴-۲) نویسه‌گردانی بر اساس تلفظ آوایی	۳۴
شکل (۵-۲) نویسه‌گردانی براساس املاء	۳۵
شکل (۶-۲) ترکیبی از سیستم نویسه‌گردانی براساس تلفظ (M_p) و براساس املاء (M_s)	۳۹
شکل (۷-۲) روش ترکیبی نویسه‌گردانی با ترکیب چند سیستم مجزا (M_i)	۴۱
شکل (۸-۲) شمای کلی سیستم نویسه‌گردانی استخراج ترجمه	۴۴
شکل (۱-۳) روند تغییرات دقت براساس حد آستانه	۵۳
شکل (۲-۳) روند تغییرات F1 براساس حد آستانه	۵۴
شکل (۳-۳) مراحل الگوریتم پیشنهادی نویسه‌گردانی	۵۶
شکل (۴-۳) نمونه‌ای از گراف همسایگی اسامی افراد	۶۰
شکل (۵-۳) تابع نگاشت فارسی به انگلیسی	۶۱
شکل (۶-۳) شبه کد الگوریتم مرحله برخط با استفاده از روش ترکیبی	۶۴

شکل ۳-۷) مقایسه روش‌های مرحله برخط..... ۶۸

شکل ۳-۸) مقایسه روش پیشنهادی Hybrid و روش Spelling..... ۶۹

شکل ۳-۹) روند صعودی دقت کلمه براساس صفحات وب در مرحله برون خط..... ۷۰

فهرست جدول‌ها

عنوان	شماره صفحه
جدول ۱-۳ (ارزیابی با حد آستانه‌های متفاوت.....	۵۲.....
جدول ۲-۳ (وزن بین کاندیدها (تعداد هر کلمه).....	۶۵.....
جدول ۳-۳ (مقایسه دقت کلمه (/.) در روشهای مرحله برخط.....	۶۷.....
جدول ۴-۳ (مقایسه دقت کلمه (/.) روش پیشنهادی Hybrid و روش Spelling.....	۶۸.....
جدول ۱-آ (پرس‌وجوهای نام.....	۷۵.....
جدول ۱-ب (پرس‌وجوهای غیرنام.....	۷۶.....
جدول ۱-ج (نمونه‌ای از پیکره اسامی فارسی به انگلیسی.....	۷۷.....

فصل ۱. مقدمه

۱-۱ مقدمه

بازیابی اطلاعات بین زبانی یکی از شاخه‌های پرکاربرد در حوزه هوش مصنوعی می‌باشد که با ترجمه پرس‌وجو و اسناد موجود در وب، دسترسی تمامی افراد به اطلاعات با زبان‌های مختلف را میسر می‌سازد. در مقوله بازیابی اطلاعات بین زبانی، به دلیل تعدد قابل توجه اسامی افراد در پرس‌وجوها، این نوع از اسامی از اهمیت ویژه‌ای برخوردار هستند. بنابراین با تشخیص اسم افراد در پرس‌وجو، صفحات شخصی افراد که شامل این نام می‌باشند، در رتبه بالاتری قرار می‌گیرند [SWZ+08]. همچنین با توجه به اینکه افراد زیادی صفحات شخصی خود را به زبان دیگری طراحی می‌کنند، می‌توان با نویسه‌گردانی آوایی^۱ پرس‌وجوها، این صفحات را در نتایج پرس‌وجو لحاظ کرد و بازیابی اطلاعات بین زبانی را انجام داد.

در این بخش، در ابتدا در مورد بازیابی اطلاعات بحث خواهد شد. پس از آن مفهوم بازیابی اطلاعات بین زبانی و کاربردهای آن و سپس چالش‌هایی که این سیستم‌ها با آن روبرو هستند، شرح داده می‌شود. همچنین مختصری در مورد تشخیص اسامی خاص در پرس‌وجوی کاربران اشاره شده و نیز نویسه‌گردانی کلمات شرح داده می‌شود. در نهایت نیز معیارهای ارزیابی سیستم‌های تشخیص نام و نویسه‌گردانی معرفی می‌گردند.

۲-۱ موتور جستجو و بازیابی اطلاعات

وب جهان‌گستر^۲ به علت توزیع‌شدگی و هزینه پایین تولید محتوا، با چالش‌های جدیدی از جمله حجم زیاد اطلاعات، رشد نمایی، پویایی زیاد، ناهمگنی و غیرساختاریافته بودن اطلاعات (شامل صوت و تصویر و غیره)، مواجه شده است [زارع ۸۸]. برای مثال حجم اطلاعات نمایه‌سازی شده در سال ۱۹۹۴ از ۱۱۰,۰۰۰ صفحه [McB94] به ۱۱/۵ میلیارد صفحه در سال

^۱ Transliteration

^۲ World Wide Web (WWW)

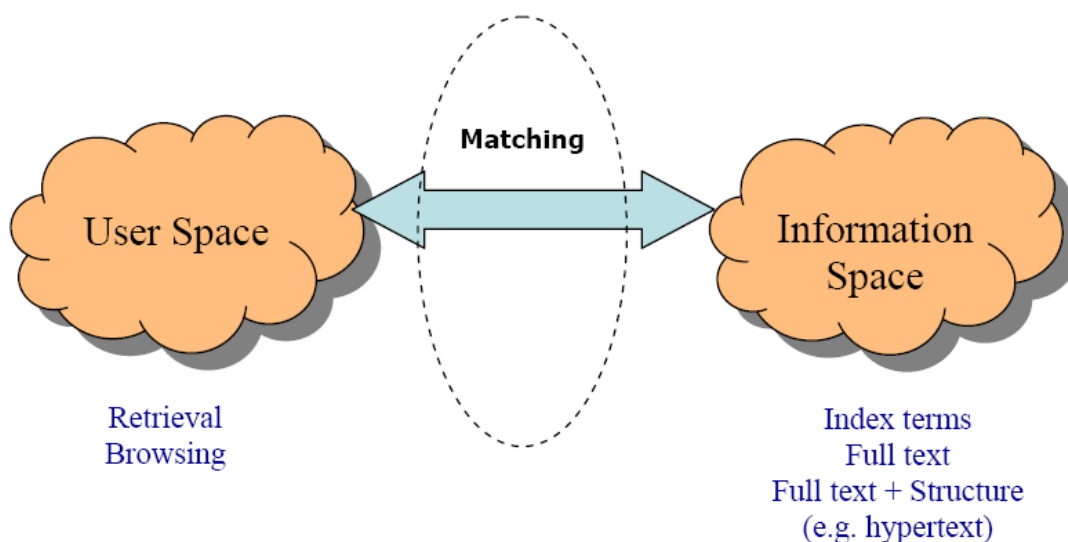
۲۰۰۵ [GS05] و پس از آن به ۲۰ میلیارد صفحه در سال ۲۰۱۱ [Web11] رشد یافته است. در حال حاضر، بهترین ابزار کارآمد برای مدیریت، بازیابی و استخراج اطلاعات مهم از این مجموعه عظیم داده‌ای، موتورهای جستجو می‌باشند. لذا حدود ۸۰٪ از کاربران اینترنت سایت‌های جدید را از طریق موتور جستجو بازیابی و استخراج می‌کنند [Inews].

۱-۲-۱ بازیابی اطلاعات^۱

براساس [زارع ۸۸]، بازیابی اطلاعات شامل استانداردها و پروتکل‌های نمایش، ذخیره‌سازی، سازماندهی و دسترسی به اقلام اطلاعاتی با هدف بازیابی کلیه اسنادی که با پرس‌وجوی کاربر مرتبط است، می‌باشد. مطابق شکل ۱-۱، بازیابی اطلاعات به دو حوزه کاربر و اطلاعات و یک حوزه میانی بنام حوزه بازیابی تقسیم می‌شود. در حوزه کاربر، نیاز اطلاعاتی کاربر بیان می‌شود که در اکثر موارد باید با توجه به زبان سیستم بازیابی این نیاز بیان شود یا این که سیستم با توجه به رفتار کاربر، تاریخچه‌ی رفتار وی و یا اطلاعات مستقیمی که از تخصص یا موارد شخصی وی دارد، او را مدل کرده و نیاز اطلاعاتی او را پیش‌بینی کند. در حوزه اطلاعات باید اطلاعات و دانش نهفته در اسناد و داده‌ها مدل‌سازی و سازماندهی شوند. حوزه بازیابی نیز فصل مشترک این دو و تطبیق دهنده نیاز اطلاعاتی کاربر با اسناد اطلاعاتی است. نیازهای اطلاعاتی یا به صورت پرس‌وجوهای برخط^۲ به حوزه بازیابی ارسال می‌شوند که در این حالت از مجموعه‌ی اسناد موجود، موارد مرتبط بازیابی می‌شود و یا این که پرس‌وجویی وجود دارد که در معرض جریان اطلاعات (مثلاً اخبار) قرار می‌گیرد و اسناد مرتبط را فیلتر و جدا می‌کند.

¹ Information Retrieval

² Online



شکل ۱-۱) نمای کلی سیستم بازیابی [BYRN99].

در فضای کاربر برای دستیابی به اطلاعات مفید، دو سیاست جستجو و مرور^۱ وجود دارد. روش جستجو برای حالتی که کاربر دقیقاً هدف خود را می‌داند، مفید می‌باشد. روش مرور وقتی مناسب است که کاربر با محتوای مورد نظر خود نا آشنا باشد. دو روش جستجو و مرور، مکمل یکدیگر بوده و تاثیر عمده‌ی خود را در حالتی که با هم استفاده شوند، خواهند داشت [زارع ۸۸].

۳-۱ بازیابی اطلاعات بین زبانی

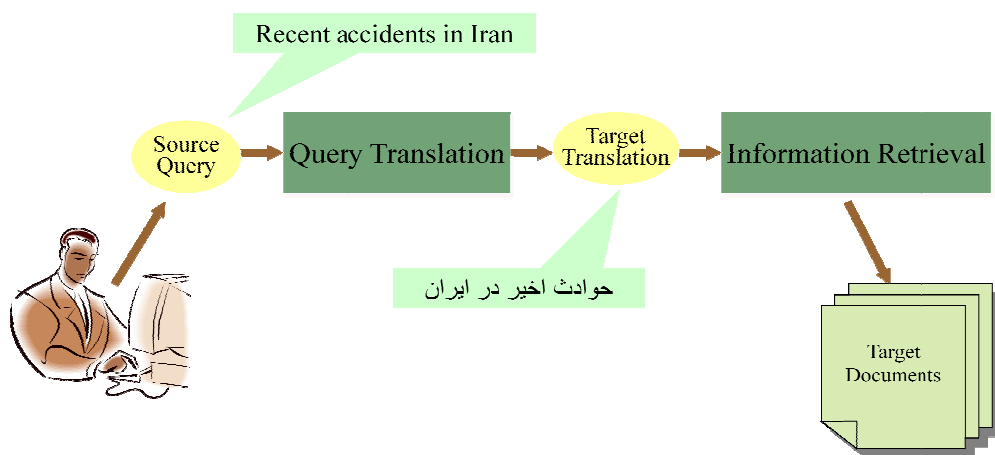
بازیابی اطلاعات بین زبانی، به نوعی از بازیابی اطلاعات گفته می‌شود که نتایج بدست آمده از پرس‌وجوی کاربر در مقایسه با پرس‌وجو دارای زبان متفاوتی است. مطابق شکل ۱-۲ هدف اصلی در بازیابی اطلاعات بین زبانی، پیدا کردن اطلاعات در زبان هدف^۲ (همانند فارسی) در پاسخ به پرس‌وجوهایی به زبان اصلی^۳ (همانند انگلیسی) می‌باشد. امروزه با رشد وسیع فناوری وب و

^۱ Browsing

^۲ Target language

^۳ Source language

افزایش روزافزون حجم اطلاعات الکترونیکی با زبان‌های گوناگون، بازیابی اطلاعات مستقل از زبان اسناد، اهمیت ویژه‌ای پیدا کرده است.



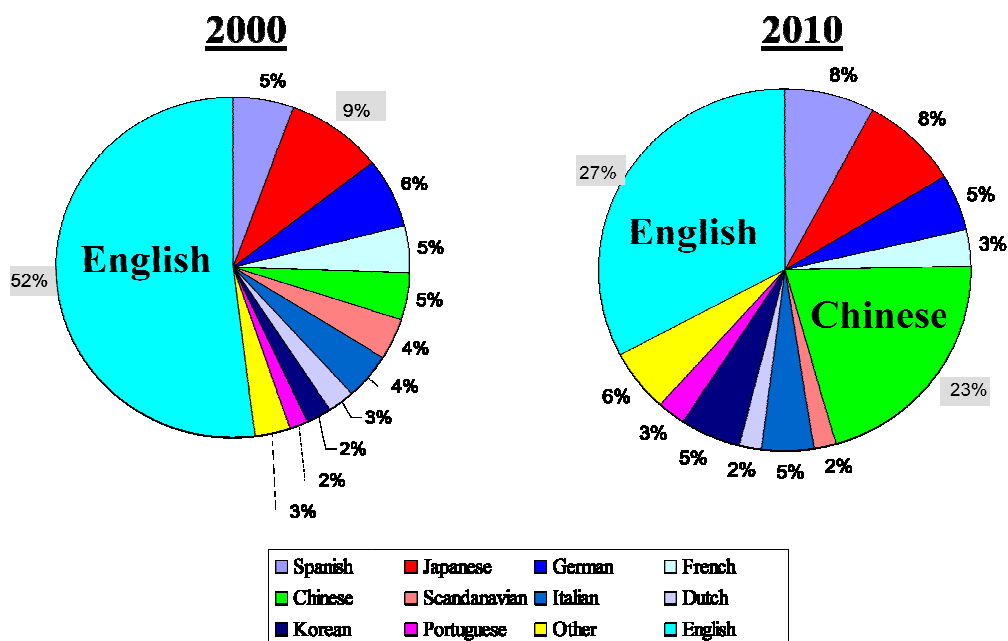
شکل ۱-۲) بازیابی اطلاعات بین زبانی.

در شکل ۱-۳، اهمیت روزافزون بازیابی اطلاعات بین زبانی مشهود است. همان‌طور که دیده می‌شود، یکی از مهم‌ترین دلایل آن محیط چندزبانه بودن^۱ وب می‌باشد که در سال‌های اخیر، به بستری برای تمامی افراد با زبان‌های مختلف تبدیل شده است.

بازیابی اطلاعات بین زبانی به منظور تبادل و اشتراک اطلاعات، به کار می‌رود. این موضوع، به عنوان نمونه در امنیت ملی^۲ و همچنین در دسترسی به اطلاعات مربوط به بیماری‌ها، نقش خود را بهتر نشان می‌دهد [RGB+03].

¹ Multilingual environment

² National security



شکل (۳-۱) چندزبانه بودن محیط وب.

این نوع سیستم‌ها زمانی کاربرد دارد که کاربرانی که توانایی خواندن چندین زبان را دارند، می‌توانند از زبانی برای نوشتن پرس‌وجویشان استفاده نمایند که برای آنها روان‌تر باشد. همچنین می‌توانند بفهمند که آیا اسنادی مرتبط به پرس‌وجوی آنها به زبان‌های دیگر وجود دارد یا خیر. از دیگر کاربردها می‌توان به بازیابی تصویر و یا ویدئو از سایت‌ها با زبان‌های گوناگون اشاره نمود [Ord97].

بازیابی اطلاعات بین زبانی شامل مقوله‌هایی از قبیل بازیابی اطلاعات، پردازش زبان طبیعی^۱، ترجمه ماشینی و خلاصه‌سازی^۲، پردازش گفتار^۳، ارتباط متقابل رایانه و انسان^۴ می‌باشد.

² Natural language processing

³ Machine translation and summarization

⁴ Speech processing

⁵ Human-Computer interaction

۱-۳-۱ چالش‌های بازیابی اطلاعات بین زبانی

همان‌طور که گفته شد، هدف از بازیابی اطلاعات بین زبانی پیدا کردن اسناد مرتبط به پرس-وجو در زبانی متفاوت با آن می‌باشد. بنابراین در نیل به این هدف، ممکن است با چالش‌های گوناگونی روبرو شود. از جمله این چالش‌ها می‌توان به ابهام کلمات در زبان‌های مختلف اشاره کرد. یک کلمه در هر زبان ممکن است دارای معانی متفاوتی باشد [هاشمی ۸۹]. برای مثال کلمه python دارای سه مفهوم کاملاً متفاوت ازدها، زبان برنامه‌نویسی و غیب‌گو می‌باشد. بنابراین با داشتن چنین ابهامی در پرس‌وجوی کاربر، با مشکل مواجه می‌شویم.

از دیگر چالش‌ها می‌توان به ترجمه عبارات و اصطلاحات^۱ اشاره نمود. هر زبان دارای اصطلاحات و عباراتی است که نمی‌توان با ترجمه تک‌تک کلمات آن به مفهوم آن پی برد [هاشمی ۸۹]. بنابراین برای داشتن دقت بالا در یک سیستم بازیابی اطلاعات بین زبانی، می‌بایست این‌گونه اصطلاحات به درستی شناخته و ترجمه شوند.

مشکل عمده دیگری که مطرح می‌گردد، زمانی است که از واژه‌نامه دوزبانه^۲ به منظور ترجمه استفاده می‌شود. در این حالت کلماتی وجود دارند که در واژه‌نامه موجود نیستند که به آنها کلمات خارج از واژه‌نامه^۳ گویند [DL05]. بنابراین سیستم بازیابی اطلاعات بین زبانی برای ترجمه این‌گونه کلمات می‌بایست از روش‌های دیگری به غیر از واژه‌نامه استفاده نماید. رفع این چالش در زمینه اسامی افراد در پرس‌وجو از جمله اهداف این پایان‌نامه می‌باشد.

همچنین از دیگر مشکلاتی که در استفاده از واژه‌نامه‌ها برای ترجمه وجود دارد، عدم وجود صرف‌های گوناگون کلمات در واژه‌نامه می‌باشد [هاشمی ۸۹]. به بیان دیگر، شکل پایه‌ای هر کلمه در واژه‌نامه قرار داده شده است و صرف‌ها و مشتقات گوناگون آن موجود نمی‌باشد.

¹ phrase

² Bilingual dictionary

³ Out of dictionary

پردازش‌های ریخت‌شناسی^۱ روی کلمات در بسیاری از زبان‌ها مانند عربی مشکلاتی را به همراه دارد [CLI05]. برای مثال حذف علائم تفکیک‌کننده^۲ در این زبان ابهاماتی را به وجود می‌آورد.

۴-۱ تشخیص نام در پرس‌وجو

با گسترش اطلاعات موجود در وب نیاز به یک موتور جستجوی قوی به منظور بازیابی این اطلاعات، احساس می‌شود. یکی از مواردی که به صحت و دقت موتور جستجو کمک می‌کند، شناخت انواع پرس‌وجوها می‌باشد.

اسامی افراد یکی از پرکاربردترین پرس‌وجوها در موتور جستجو می‌باشند. براساس [SWZ+08]، با استفاده از داده‌های جمع‌آوری شده از یک موتور جستجوی تجاری، ۴~۲٪ از پرس‌وجوهای کاربران منحصراً اسامی افراد بوده و ۳۰٪ از آنها، پرس‌وجوهایی هستند که شامل اسامی افراد و کلمات دیگری غیر از این اسامی می‌باشند [AVG05]. این مقادیر در حالتی که کاربران نتایج حاصل از یک موتور جستجو را نامناسب ارزیابی کرده و موتور جستجوی دیگری را انتخاب می‌نمایند، افزایش می‌یابد. تعداد زیاد پرس‌وجوهای اسامی افراد از آنجا ناشی می‌شود که اکثر افراد، نام خود را به عنوان پرس‌وجو وارد کرده و عملکرد موتور جستجو را از طریق نتایج حاصل، ارزیابی می‌کنند. در صورت شناخت موفقیت‌آمیز این نوع از پرس‌وجوها، موتور جستجو می‌تواند عملکرد بهتری در رتبه‌بندی نتایج داشته باشد. به طور کلی، اگر بتوان پرس‌وجوی دارای نام را تشخیص داد، می‌توان خصیصه‌هایی هم‌چون صفحه شخصی^۳ افراد و یا تعداد کلمات پرس-وجو در صفحه، را در رتبه‌بندی نتایج اعمال کرد [SWZ+08]. همچنین با توجه به این که تعداد زیادی از کاربران صفحات شخصی خود را به زبان‌های دیگر طراحی می‌کنند، می‌توان پس از

¹ Morphological analyzing

² Punctuation mark

³ Personal homepage