

نشریه پژوهش‌های زبان‌شناسی

سال هفتم، شماره دوم، شماره ترتیبی ۱۳، پاییز و زمستان ۱۳۹۴

تاریخ وصول: ۱۳۹۳/۱۱/۱

تاریخ اصلاحات: ۱۳۹۴/۶/۱۶

تاریخ پذیرش: ۱۳۹۵/۱/۳۰

صص ۱۴-۱

استخراج خودکار معادل‌های واژگانی از پیکره‌های دو زبانه موازی

طیبه موسوی میانگاه*

مهشید شکیبیا**

چکیده

امروزه به مدد ظهور انواع فناوری‌های رایانه‌ای، پیکره‌های زبانی نقش بسیار مهمی در حل انواع مختلف مسائل زبانشناختی ایفا می‌کنند. پیکره‌های دو زبانه موازی در سطح جمله و در سطح واژه می‌توانند برای بازیابی واحدهای تک‌واژه‌ای و یا حتی چند واژه‌ای براحتی مورد استفاده قرار گیرند که این امر کاربردهای مفیدی در حوزه‌های مختلف رایانه و زبان خواهد داشت. هدف این مقاله به کارگیری یک پیکره موازی انگلیسی^۱ فارسی از قبل طراحی شده در جهت ساخت یک مطابقه^۱ (کشف اللغات) دو زبانه کارآمد با استفاده از آماره اطلاعات متقابل است. در اینجا از آماره اطلاعات متقابل استفاده می‌شود تا همترازی در سطح واژه بین جملات انگلیسی و فارسی پیکره مورد نظر صورت گیرد. یک پیکره زبانی همتراز شده در سطح واژه مسلماً کاربردهای زیادی از جمله در تهیه نرم‌افزار حافظه ترجمه، مدیریت مجموعه اصطلاحات، بازیابی اطلاعات دوزبانه، سیستم ترجمه ماشینی مبتنی بر آمار و مانند آن دارد. با استفاده از یک الگوریتم ابتکاری آزمایشی ترتیب داده شده و مقایسه‌ای بین برونداد همترازسازی خودکار با جملات همتراز شده توسط مترجم انسانی صورت گرفت. نتایج این آزمایش نشان داد که برنامه مطابقه گزارش شده در این تحقیق می‌تواند صحتی معادل ۷۵ درصد را به دست آورد.

کلیدواژه‌ها: مطابقه دو زبانه، ترجمه خودکار، بازیافت اطلاعات، استخراج معادل‌های واژگانی، پیکره موازی انگلیسی-فارسی

¹ concordance

Mosavit@pnu.ac.ir

Mahshid.Shakiba@gmail.com

* دانشیار گروه زبان‌شناسی دانشگاه پیام نور

** کارشناس ارشد گروه کامپیوتر دانشگاه پیام نور

۱. مقدمه

در سالهای اخیر ساخت و استفاده از انواع مختلف پیکره‌های زبانی، از جمله کاربردهای رایانه‌ای هستند که در اختیار پژوهشگران حوزه ترجمه ماشینی، آموزش و یادگیری زبان، مطالعات ترجمه، بازیافت اطلاعات دو زبانه، واژه‌نگاری و غیره قرار گرفته‌اند. تا قبل از توسعه فناوری اطلاعات و ارتباطات، پیکره‌ها و نرم‌افزارهای کشف اللغات بندرت در دسترس پژوهشگران زبان قرار داشتند تا بتوانند در مورد زبان، محتوا، و یا امور مربوط به ترجمه اطلاعاتی را از درون آن‌ها استخراج نمایند.

ولی امروزه با ظهور فناوری‌های رایانه‌ای و تولید حجم وسیع متون ترجمه، منابع پیکره‌ای قابل دسترس برای محققان اشتیاق زیادی را در آنان برای ساخت و بهره‌برداری از پیکره‌ها به وجود آورده است. پیکره‌های ترجمه یا پیکره‌های موازی در رشته زبان‌شناسی رایانه‌ای بویژه در سیستم‌های ترجمه ماشینی کاربردهای وسیعی دارند. در حقیقت، پیکره‌های موازی ابزار اصلی برای انتخاب ترجمه‌های ممکن واحدهای واژگانی و متعاقب آن یافتن محتمل‌ترین ترجمه‌ها در یک سیستم ترجمه ماشینی مبتنی بر آمار هستند.

مطابق نرم‌افزاری است برای مشاهده پیکره که فهرستی از رخدادهای یک واژه خاص، قسمتی از یک واژه یا ترکیبی از چند واژه را در بافت نشان می‌دهد. به واحد جستار کلیدواژه گفته می‌شود. رایج‌ترین شکل نمایش در یک مطابقه نشان‌دادن تعدادی سطر است با ساختار کلیدواژه در بافت که در آن هر کلیدواژه با طول بافت معینی از هر دو طرف در مرکز سطر نشان داده می‌شود. نمونه‌ای از یک مطابقه دو زبانه انگلیسی^۱ فارسی (موسوی میانگانه، ۲۰۰۹) در شکل (۱) نشان داده شده است.

The screenshot shows the Machine Translation V 1.1.0 interface. At the top, there's a search bar with the text "حقوق بشر" (Human Rights) and a search button. Below it, there's a result list. The main part of the screen is a table with two columns: English and Persian. The English text is a message from the UN Secretary-General on the International Human Rights Day. The Persian text is a translation of this message. The table has a header row with "English" and "Persian" and a "Number" column. The table contains 77 rows of text. The Persian text is a translation of the English text, with some words in brackets indicating corrections or suggestions. The interface also includes a "Delete Selected" button and a "Number" field showing 77.

شکل (۱) رکوردهای تولیدشده توسط مطابقه موازی در سطح جمله برای جستار "حقوق بشر"

¹ Keyword in Context

تمام پیکره‌های ماشین خوان^۱ به محققان در محاسبه فراوانی رخداد واژه یا عبارت مورد جستجو (جستار) کمک می‌کنند تا با انتخاب واژه یا عبارت مورد نظر تمام رکوردها (سطور) مطابقتی که در آن جستار مورد نظر وجود دارد در صفحه نمایش ظاهر شود. بعلاوه، اکثر بسته‌های تحلیل پیکره شامل مطابقتی نیز می‌باشند که محققان را در یافتن تمام رخداد‌های جستارهای مورد نظر با امکان مرتب کردن داده‌های نمایش داده شده در صفحه و نیز همراه با محدوده بافت متعلقه از سمت چپ و راست یاری می‌نمایند (هیرویوکی و توشیکو، ۱۹۹۶). واژه‌های خاص، قسمتی از یک واژه، یا گروهی از واژه‌ها که توسط برنامه مطابقتی نمایش داده می‌شوند همه از یک پیکره متنی استخراج می‌شوند. بدیهی است که هر چه حجم پیکره بیشتر باشد، برونداد برنامه دقیق‌تر خواهد بود. در واقع، مطابقتی‌های دوزبانه از جمله ابزاری هستند که کار سخت ترجمه را آسان‌تر می‌نمایند.

یک مطابقتی موازی از این جهت با یک مطابقتی یک‌زبانه متفاوت است که در مورد اول برنامه مطابقتی براساس مجموعه‌ای از متون یک‌زبانه ساخته می‌شود مثل برنامه مایکروکانکورد (جونز، ۱۹۸۶) و ورد اسمیت (اسکات، ۲۰۰۰)، در حالی که در مورد دوم برنامه مطابقتی براساس مجموعه‌ای از متون موازی دو یا چندزبانه ساخته می‌شود. این بدان معناست که در مطابقتی‌های موازی کاربر می‌تواند جستار خود را به یک زبان نوشته و جملات متناظر برای آن جستار را به زبان دیگر همان زبان بلکه به زبان یا زبان‌های دیگر دریافت نماید. دو نمونه از پرکاربردترین مطابقتی‌های موازی پاراکانکورد (بارلو، ۲۰۰۲) و مولتی کانکورد (جونز، ۱۹۹۸) هستند.

از طرف دیگر، مطابقتی موازی در سطح واژه نوع خاصی از مطابقتی دوزبانه است که در آن جستجوی کاربر برای یک زنجیره در زبان مبدأ منتهی به یافتن جملات متناظر با آن جستار در هر دو زبان مبدأ و مقصد شده، در حالی که زنجیره جستار در جمله زبان مبدأ و ترجمه آن در جمله مقصد پررنگ و کاملاً مشخص شده است. داشتن این توانایی برای مطابقتی‌ها مستلزم دستیابی به فنون پیشرفته همتراسازی در سطح واژه است. البته در بسیاری از موارد هم شناسایی ترجمه دقیق جستار مورد نظر امکان پذیر نمی‌باشد.

در سال‌های اخیر بسیاری از سیستم‌های ترجمه ماشینی با استفاده از همتراسازی در سطح واژه، پیکره‌های موازی که توسط مدل‌های آی‌بی‌ام (براون و همکاران، ۱۹۹۳) و بسته جیزاپلاس‌پلاس^۲ (اوج، ۲۰۰۰ و اوج و نی، ۲۰۰۳) تولید شده طراحی شده‌اند. اما این پژوهش روشی نسبتاً ابتکاری و نوین برای استخراج تناظرها در سطح واژه از پیکره موازی همتراسازی شده در سطح جمله و با استفاده از آماره اطلاعات متقابل ارائه می‌نماید. در بخش‌های زیر پس از مرور تحقیقات انجام گرفته در این حوزه، پیکره دوزبانه‌ای که برای این پژوهش مورد استفاده قرار گرفته معرفی می‌شود و متعاقب آن روش تحقیق بکاررفته برای انجام آزمایشی در رابطه با ارزیابی کیفی این برنامه مطابقتی همراه با ارائه الگوریتم آن مورد بحث قرار خواهد گرفت. آزمایش و نتایج حاصل از آن در بخش ۴ و ۵ گزارش خواهند شد. این مقاله با ملاحظات جمع‌بندی و بحث در مورد کارهای آتی به پایان خواهد رسید.

^۱ Machine readable corpora

^۲ aligning

^۳ GIZ/A++ Toolkit

^۴ Mutual information

۲. مروری بر پژوهش‌های پیشین

تاکنون تحقیقات زیادی در زمینه پیکره‌های موازی صورت گرفته‌است. تلاش تمام این تحقیقات به سمت استفاده از کاربردهای چنین پیکره‌هایی به منظور حل مسائل مختلف زبانشناسی رایانشی مانند ترجمه ماشینی مبتنی بر آمار (براون و همکاران، ۱۹۹۳)، بازیابی اطلاعات دوزبانه (نظارات و موسوی میانگانه، ۱۳۹۰)، یادگیری زبان (وو و همکاران، ۲۰۰۳)، ترجمه انسانی (موسوی میانگانه، ۲۰۰۶) و نظیر آن بوده است. برخی از محققان سعی کرده‌اند تا روش‌های جدیدی برای همتراسازی پیکره‌های دوزبانه در سطح جمله بیابند (سیمارد و پلامندن، ۱۹۹۸ و اوچ و نی، ۲۰۰۱). برخی دیگر از آنها به سمت استفاده از روش‌های آماری (گیل و چرچ، ۱۹۹۱؛ کاپیک، ۱۹۹۳؛ داگان و همکاران، ۱۹۹۳؛ اینو و ناگایتو، ۱۹۹۳ و فانگ، ۱۹۹۵) یا زبانی (یاماماتو و ساکاماتو، ۱۹۹۳؛ کومانو و هیراکاوا، ۱۹۹۴ و ایشیماتو و ناگاوا، ۱۹۹۴) برای همتراسازی پیکره‌های دوزبانه در سطح واژه روی آورده‌اند. روش‌های آماری از فراوانی واژه‌ها استفاده می‌کنند و میزان ارتباط بین واژه‌های متناظر در دوزبان موجود در پیکره موازی را محاسبه می‌نمایند. این در حالیست که روش‌های زبانی تناظرهای واژگانی در دو زبان را با استفاده از یک واژه‌نامه دوزبانه پیدا می‌کنند.

ونگ لکسم یک برنامه مطابقه موازی انگلیسی - چینی طراحی نموده است و کاربردهای آموزشی آن در زبان‌های انگلیسی و چینی را از طریق مثال‌هایی از آزمایش‌های آموزش و یادگیری نشان داده است، و بدین ترتیب روش یادگیری داده‌محور را پیاده‌سازی نموده است. در واقع، در پژوهش او، این ارزش آموزشی مطابقه موازی است که مورد توجه قرار گرفته و نه روش یافتن تناظرهای واژگانی در پیکره موازی (ونگ، ۲۰۰۱).

همینطور یک سیستم مطابقه انگلیسی-چینی مبتنی بر شبکه بنام توتال‌ریکال^۱ به منظور افزایش استفاده مجدد از ترجمه‌ها و نیز تشویق به استفاده از منابع موثق و طبیعی در نگارش زبان دوم توسط وو و همکارانش طراحی شده است. توتال‌ریکال پیشرفته‌تر از مطابقه‌های قبلی است چرا که در این سیستم کاربر نه تنها می‌تواند تمام سوابق مربوط به واحد واژگانی مورد جستجو را ببیند بلکه می‌تواند معادل‌های ترجمانی آن را به صورت برجسته مشاهده نماید. این سیستم در حقیقت نوعی مطابقه موازی همتراساز شده در سطح واژه است که نسبت به مطابقه‌های قبلی کاربردهای بسیار بیشتری دارد (وو و همکاران، ۲۰۰۳).

پاراکانکورد مطابقه موازی و چندزبانه‌ای است که چهار متن موازی از چهار زبان متفاوت، یا یک متن اصلی و سه ترجمه متفاوت از سه زبان دیگر را می‌پذیرد. این برنامه دربرگیرنده نوار ابزاری برای برجسته‌سازی ترجمه‌های کاندید شده و نیز یک بخش خودکار بنام هات وردز است که از اطلاعات بسامدی استفاده می‌کند تا اطلاعاتی در مورد ترجمه‌های ممکن جستار موردنظر فراهم نماید (بارلو، ۲۰۰۲).

موزر و همکارانش یک سیستم همتراسازی در سطح واژه برای ترجمه ماشینی آماری ارائه دادند که بطور همزمان جمله متن مقصد را طوری مجدداً مرتب‌سازی می‌کند که با ترتیب واژه‌ها در جمله متن مقصد متناظر خود منطبق باشد. ایده اصلی کار آن‌ها تولید یک الگوریتم یکنواخت بین جمله مقصد و یک شبکه‌بندی تبدیلی بود که نشان‌دهنده مرتب‌سازی‌های مجدد متفاوت واژه‌ها در جمله مبدأ باشد. آن‌ها نشان دادند که سیستم‌های ترجمه‌ای که با روش

¹ TotalRecall

² hot words

³ monotonic

پیشنهادی آن‌ها کار می‌کنند بهتر و یا حداقل در همان حد سیستم‌هایی که با بسته نرم افزاری جیزاپلاس پلاس کار می‌کنند عمل می‌کنند (مازر و همکاران، ۲۰۰۶).

۳. پیکره موازی انگلیسی-فارسی

پیکره موازی انگلیسی-فارسی در ابتدا به صورت بانک داده‌ای متنی متشکل از متون اصلی به زبان انگلیسی و ترجمه‌های آن‌ها به زبان فارسی و نیز متون اصلی به زبان فارسی و ترجمه‌های آن‌ها به زبان انگلیسی گردآوری شد. اگرچه میزان دسترسی به متون دوزبانه که شامل زبان فارسی باشند به دلیل تراکم پایین متون فارسی در سراسر دنیا و وجود نداشتن چنین متونی در برخی گونه‌ها و زمینه‌های خاص بسیار پایین می‌باشد، ما موفق به تهیه بانک داده‌ای دوزبانه نسبتاً حجیمی متشکل از صدهزار جمله انگلیسی و فارسی شدیم. متون در این پیکره به انواع زیر طبقه‌بندی شده‌اند: مذهبی، ادبی، سیاسی، اقتصادی، علمی، شعر، اصطلاحات و ضرب‌المثلهای، ورزشی، متفرقه، پزشکی، و فرهنگی. نوع متن با جستجوی هر جمله در کنار آن ظاهر می‌شود. پایگاه داده‌ای این پیکره به صورت اکسس^۱ و اس کیوال^۲ موجود می‌باشد و مطابقه مبتنی بر آن نیز به صورت تحت ویندوز و تحت وب تهیه شده است. پیکره تولید شده پیکره‌ای در حال پیشرفت است، بدین معنا که پیکره‌ایست باز که بر حسب نیاز و با گذشت زمان مطالبی به آن اضافه خواهد شد (موسوی میانگاه، ۲۰۰۹).

آماده‌سازی و همترازسازی پیکره

متون خامی که از منابع گوناگون استخراج می‌شوند برای ورود به پیکره باید پیش پردازش شوند. دانلود کردن، تبدیل فرمت و هنجار سازی (عادی سازی) متون از جمله مراحل بسیار وقت گیر در آماده‌سازی پیکره به شمار می‌رود. برخی از صفحات که نامربوط هستند و نیز تمام شکل‌ها، جداول و عکس‌ها قبل از اینکه وارد پیکره شوند باید حذف گردند. در برخی موارد که یک جمله یا قسمتی از آن ترجمه نشده باشد، قسمت‌های ترجمه نشده باید حذف شوند. بعد از بازبینی و تأیید، تمام متون به طور هماهنگ به فرمت ایکس ام ال^۳ رمز گذاری می‌شوند تا پیکره بتواند مستقل از کاربرد بوده و از طریق اینترنت به راحتی قابل مبادله باشد.

به منظور استخراج اطلاعات از پیکره موازی، لازم است متون دوزبانه این پیکره در مرحله نخست در سطوح پاراگراف، جمله و واژه همتراز شوند. منظور از همتراز سازی ایجاد وابستگی بین قطعاتی از متن در یک زبان و ترجمه‌ها یا متون معادل آن‌ها در زبان دیگر است. همتراز سازی در سطح پاراگراف کار نسبتاً ساده‌ای است، چون مرزهای پاراگرافی معمولاً مشخص هستند، در ضمن اینکه این نوع همتراز سازی برای استفاده‌های بعدی از پیکره‌های موازی در امر پژوهش خیلی هم مفید به نظر نمی‌رسند. همتراز سازی در سطح واژه به معنای تعیین کردن جفت‌واژه‌های متناظر در دو زبان می‌باشد. این نوع همتراز سازی کاری بسیار دشوار بوده و نیاز به الگوریتم‌های پیچیده‌ای دارد. بنابراین اکثر پیکره‌های موازی در سطح جمله همتراز می‌شوند.

¹ Access

² SQL

³ XML (Extendible Mark-up Language)

در پیکره موازی انگلیسی-فارسی طراحی شده توسط نویسنده، همتراسازی در سطح جملات گرچه می‌توانست با روش‌های کاملاً خودکار یا نیمه‌خودکار انجام شود (رنسیک، ۱۹۹۸، ۱۹۹۹)، لیکن تماماً به طور دستی انجام شده است تا صحت همتراسازی به ۱۰۰ درصد برسد. قصد ما این است که این پیکره در حل مسائلی مانند پردازش خودکار متن و ترجمه ماشینی آماری و ساختن نرم افزار حافظه ترجمه که در آن‌ها دقت بالا بسیار مهم می‌باشد مورد استفاده قرار گیرد. پیکره موازی انگلیسی^۱ فارسی فوق در انجمن منابع زبانی اروپا به ثبت رسیده و دستیابی به آن از طریق مراجعه به کاتالوگ این انجمن که آدرس آن در زیر آورده می‌شود امکان‌پذیر می‌باشد:

http://catalog.elra.info/product_info.php?products_id=1111

ساخت یک پیکره به‌خودی‌خود هدف محسوب نمی‌شود بلکه معمولاً به عنوان قسمتی از یک پروژه تحقیقاتی در نظر گرفته می‌شود. در واقع هنگامی که ساخت پیکره به پایان می‌رسد کار اصلی شروع می‌شود.

۴. الگوریتم همتراسازی در سطح واژه

پس از این که همتراسازی پیکره مورد نظر در سطح جمله صورت گرفت، مرحله بعدی تشخیص تناظرهای واژگانی در جفت جملات متناظر است. به عبارت دیگر، هدف اصلی در اینجا این است که مشخص شود کدام واژه در جمله انگلیسی با کدام واژه در جمله فارسی متناظر آن مطابقت دارد. این دو واژه معادل در جفت جملات متناظر برجسته یا پررنگ می‌شوند.

از آنجاکه میان طول یک جستار و سرعت پاسخ‌دهی رابطه عکس برقرار است، یعنی هرچه قدر یک جستار تعداد واژه‌های بیشتری را دربرگیرد احتمال تهی بودن نتیجه بیشتر است، در این روش عمدتاً یک واژه به عنوان جستار در نظر گرفته شده است. همچنین، از آنجایی که واحد اصلی جستار در امور مربوط به نمایه از قبیل بازیابی اطلاعات در بیشتر موارد اسم‌های مجزا هستند، تمام واژه‌های قاموسی زبان فارسی باید با تمام واژه‌های محتوایی زبان انگلیسی مقایسه شوند تا از محاسبات اضافی و اختلال ممانعت به عمل آید. تشخیص واژه‌های محتوایی از واژه‌های دستوری از طریق برچسب‌زن اجزای کلام که با استفاده از سیستم برچسب‌زن کلاز^۲ در مورد واژه‌های انگلیسی اعمال می‌شود براحتی امکان‌پذیر است. در مورد زبان فارسی نیز از سیستم برچسب‌زن اجزای کلام فارسی تگ^۳ استفاده می‌شود (موسوی میانگه، ۲۰۱۵). از آنجا که واژه‌های دستوری در هر زبانی مقولات بسته‌ای را تشکیل می‌دهند با خارج کردن این مقولات از واژگان، بقیه واژه‌ها محتوایی در نظر گرفته شده‌اند. در مواردی که بین این دو گروه واژه‌های محتوایی و دستوری هم‌پوشانی وجود داشته نیز سیستم‌های برچسب‌زن مذکور که عمدتاً سیستم‌های مبتنی بر قاعده هستند تا حد زیادی قادر به تشخیص آن‌ها بوده‌اند.

اکنون زمان آن رسیده تا مطابقه موازی در سطح جمله به شکل یک بسته نرم‌افزاری ساخته شود تا کاربر قادر باشد به وسیله آن واژه خاصی را در زبان انگلیسی یا فارسی جستجو کرده و فهرستی از تمام جملات به زبان مورد جستجو شامل

¹ content words

² function words

³ CLAWS POS tagger

⁴ FarsiTag

آن واژه خاص همراه با جملات متناظر به زبان دیگر را دریافت نماید و این در حالی است که جستار موردنظر و معادل آن به زبان دیگر به طور برجسته نشان داده شده است. به عبارت دیگر، جستجو می‌تواند براساس هر کدام از زبان‌های موجود در پیکره انجام شود.

اگر دو جمله انگلیسی و فارسی ترجمه یکدیگر باشند، انتظار می‌رود که یک واژه فرضی در جمله فارسی دارای یک ترجمه انگلیسی در جمله انگلیسی متناظر خود داشته باشد. جدول (۱) پیوندهای مناسب در دو جمله متناظر انگلیسی و فارسی را که توسط مترجم انسانی همتراز شده نشان می‌دهد.

For true colonialism to exist two conditions are necessary.

۱ ۲ ۳ ۴ ۵ ۶ ۷ ۸

برای وجود استعمار حقیقی دو شرط لازم است.

۸ ۷ ۶ ۵ ۴ ۳ ۲ ۱

جدول (۱) همترازسازی در سطح واژه در یک نمونه ساده از جفت جمله انگلیسی و فارسی

English → Persian
1 → 1
2 → 4
3 → 3
4 → 2
5 → 5
6 → 6
7 → 8
8 → 7

لازم به ذکر است که زبان فارسی زبانی هسته‌آغازین است که از راست به چپ نوشته می‌شود. در دو جمله همتراز انگلیسی^۰ فارسی، یافتن ترجمه مناسب یک واژه فرضی فارسی در میان واژه‌های انگلیسی جمله معادل آن برای یک جستجوگر غیرانسانی تقریباً غیرممکن است، چرا که ترتیب عناصر در جملات انگلیسی و فارسی با هم تطابق ندارد. از این رو، ما روشی نسبتاً بدیع برای یافتن مناسب‌ترین و محتمل‌ترین معادل هر واژه با استفاده از آماره اطلاعات متقابل ارائه داده‌ایم. اطلاعات متقابل اساساً برای محاسبه میزان همبستگی میان واژه‌ها با استفاده از آمار هم‌رخدادی واژه‌ها بکار برده می‌شود و به صورت فرمول شماره (۱) قابل تعریف است (چرچ و هنکس، ۱۹۹۰):

(۱)

$$MI(x, y) = \log_2 \frac{P(x, y)}{P(x)P(y)} = \log_2 \frac{Nf_w(x, y)}{f(x)f(y)}$$

در اینجا x و y واژه‌های فرضی در بافت هستند. احتمال‌های $p(x)$ و $p(y)$ با استفاده از شمارش تعداد رخداد‌های x و y در پیکره یعنی $f(x)$ و $f(y)$ محاسبه می‌شوند. N حجم پیکره را نشان می‌دهد. $P(x, y)$ با شمارش تعداد دفعاتی که x و y در بافت یکسانی ظاهر می‌شوند (در اینجا منظور از بافت همان رکورد است) محاسبه می‌شود. از آنجا که الگوریتم پیاده‌سازی شده برای این برنامه جمله‌ها را شناسایی نموده و واژه‌های موجود در یک جمله که جستار به آن تعلق دارد را به عنوان بافت زبانی آن جستار در نظر می‌گیرد، در این تحقیق هر رکورد شامل یک جمله انگلیسی و معادل فارسی آن است.

استفاده از مقادیر اطلاعات متقابل براساس این فرض است که هنگامی که دو واژه انگلیسی و فارسی در محدوده بافت معینی (که در اینجا یک رکورد است) با فراوانی بالایی هم‌رخداد می‌شوند، احتمال این که آن‌ها ترجمه یکدیگر باشند بیشتر می‌شود. درحقیقت، برای یافتن معادل فارسی هر واژه محتوایی در زبان انگلیسی این برنامه جملاتی را در پیکره جستجو می‌کند که در آن‌ها واژه انگلیسی مورد نظر وجود داشته باشد. سپس برنامه مقدار اطلاعات متقابل آن واژه انگلیسی با تمام واژه‌های محتوایی زبان فارسی که در جملات فارسی متناظر وجود دارند را محاسبه می‌نماید. به‌عنوان مثال، جمله هم‌تراز شده زیر را که در آن ترجمه واژه انگلیسی که زیر آن خط کشیده شده باید از جمله فارسی متناظر آن استخراج شود را در نظر بگیرید:

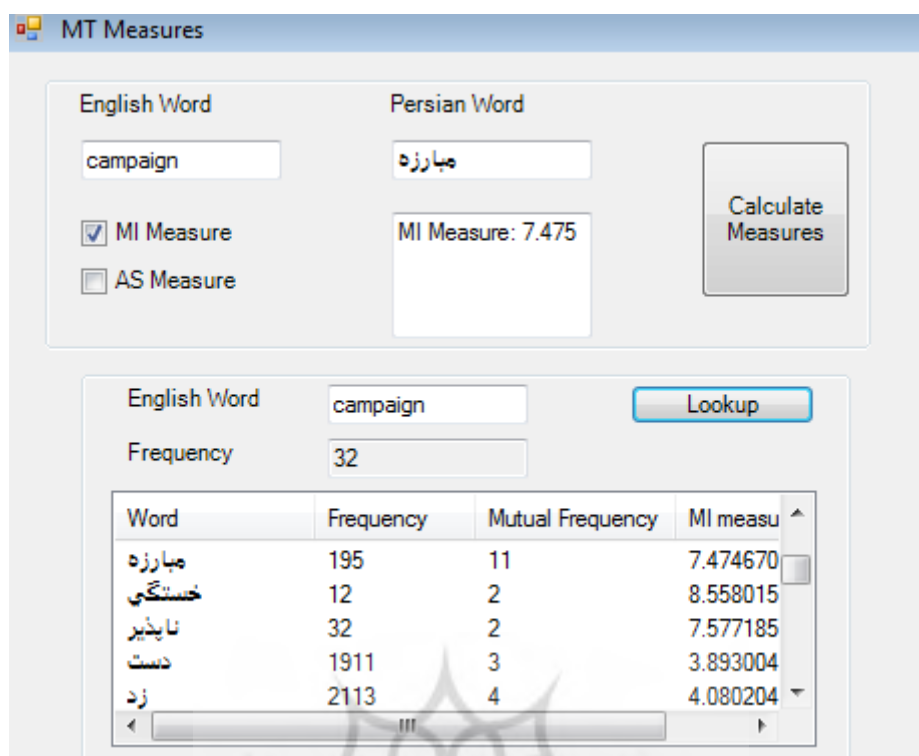
جدول (۲) نمایش یک رکورد تولیدشده توسط پیکره برای جستار "campaign"

English sentence	Persian sentence
UN has launched a global campaign against the poverty.	سازمان ملل مبارزه جهانی علیه فقر را آغاز کرده است.

در جدول (۲) مقدار اطلاعات متقابل واژه *campaign* و تمام واژه‌های محتوایی فارسی در طرف راست محاسبه شده‌اند و جفت واژه‌هایی که همبستگی زیادتری با یکدیگر دارند بعنوان ترجمه‌های یکدیگر انتخاب می‌شوند و بدین ترتیب ترجمه‌هایی که احتمال درست بودنشان کمتر است حذف می‌شوند. نتایج این محاسبات را می‌توان در جدول شماره (۳) مشاهده نمود. روش اجرای این برنامه برای واژه *campaign* در محیط نرم‌افزاری در شکل ۲ نمایش داده شده است.

جدول (۳) مقدار اطلاعات متقابل واژه "campaign" و ۸ واژه در جمله فارسی متناظر

Word x	Word y	f(x)	f(y)	f(x, y)	MI
campaign	سازمان	۳۲	۱۱۳۰	۴	۴/۷۰۶
campaign	ملل	۳۲	۱۶۰۴	۷	۴/۹۱۵
campaign	مبارزه	۳۲	۱۹۵	۱۱	۷/۴۷۴۶
campaign	جهانی	۳۲	۳۱۰	۲	۴/۴۲۸
campaign	علیه	۳۲	۱۵۹	۴	۵/۴۹
campaign	فقر	۳۲	۱۷۳	۲	۵/۸۹
campaign	آغاز	۳۲	۱۷۷	۲	۵/۸۶۶
campaign	کرده است	۳۲	۴۶۳۵	۵	۵/۹۴۳



شکل ۲) روش اجرای برنامه برای واژه campaign در محیط نرم افزاری

از جدول شماره (۳) می‌توان دریافت که برای این جمله بیشترین مقدار اطلاعات متقابل به جفت "campaign" و "مبارزه" با مقدار $7/4746$ متعلق می‌باشد و مقدار اطلاعات متقابل دیگر جفت واژه‌ها بسیار پایین‌تر از آن است. بدین ترتیب، این الگوریتم قادر خواهد بود مناسب‌ترین ترجمه یک واژه انگلیسی را در جمله فارسی متناظر آن برجسته نماید و با این کار یک مطابقت مستقل از زبان کاملاً مبتنی بر آمار تولید شده است. البته جستار می‌تواند هم به زبان انگلیسی و هم به زبان فارسی باشد. تعداد محاسبات موردنیاز برای یافتن معادل یک واژه خاص در یک جمله به یک زبان به تعداد واژه‌های محتوایی در جمله متناظر آن به زبان دیگر بستگی دارد.

۵. نتایج آزمایش

به منظور ارزیابی میزان تأثیر این الگوریتم، آزمایشی بر روی پیکره موازی موجود ترتیب داده شد. در این آزمایش تنها واژه‌های محتوایی در نظر گرفته شدند چراکه واژه‌های اصلی که توسط کاربران در تقریباً تمام مطابقت‌ها مورد جستجو قرار می‌گیرند همین نوع واژه‌ها هستند. هرچند واژه‌های دستوری نیز به انسجام متن کمک می‌کنند، اما همانطور که می‌دانیم در موتورهای جستجو و به‌طور کلی در ارزیابی اطلاعات چنین واژه‌هایی قابل ملاحظه نیستند و اکثر قریب به اتفاق جستارها در زبان‌های مختلف عمدتاً عباراتی متشکل از واژه‌های محتوایی مانند اسم و صفت هستند.

پیکره آزمون استفاده‌شده برای ارزیابی عملکرد این آزمایش مبتنی بر الگوی پیشنهادی شامل مجموعه‌ای از ۱۰۰ واژه انگلیسی در نقش جستار است که به‌عنوان درونداد به الگوریتم داده می‌شود و برونداد سیستم مطابقت در سطح جمله است که در آن کاربر قادر است مجموعه‌ای از جفت جملات به انگلیسی و فارسی را همراه با جستارهای مورد نظر و

ترجمه‌هایشان که برجسته‌نمایش داده می‌شوند را مشاهده نمایید. با هر واژه انگلیسی در پیکره‌آزمون فهرستی از ترجمه‌های ممکن ارائه می‌شود و انتخاب محتمل‌ترین آن‌ها هدف نهایی سیستم است.

در سیستم‌های بازیابی اطلاعات، معیار دقت و بازخوانی و معیارهایی شبیه به آن‌ها به عنوان معیارهای اصلی ارزیابی به کار می‌روند: معیار دقت به حاصل تقسیم «تعداد مستندات بازیابی شده واقعاً با ربط» بر «تعداد کل مستندات بازیابی شده» گفته می‌شود. معیار بازخوانی به حاصل تقسیم «تعداد مستندات بازیابی شده با مرتبط بر تعداد مستندات مرتبط موجود در مجموعه اطلاعاتی» گفته می‌شود. در واقع دقت درصد متون بازیابی‌شده مرتبط و فراخوانی درصد ارتباط متون بازیابی شده را نشان می‌دهند.

در این آزمایش عملکرد برنامه‌ای که سعی دارد واحدهای واژگانی زبان انگلیسی را با استفاده از پیکره ترجمه نماید براساس دو نوع معیار ارزشیابی، یعنی دقت و فراخوانی، طبق فرمول‌های (۲) و (۳) مورد ارزیابی قرار گرفت.

(۲)

$$Precision = \frac{\text{Number of correctly translated strings}}{\text{Total number of strings translated by program}}$$

(۳)

$$Recall = \frac{\text{Number of correctly translated strings}}{\text{Total number of strings in the test set}}$$

دقت میزان دقیق بودن یا درستی را نشان می‌دهد، درحالی که فراخوانی میزان کامل بودن را نشان می‌دهد. از آنجایی که در این آزمایش پاسخ‌های تولیدشده (تعداد زنجیره‌های تولیدشده توسط برنامه) با کل پاسخ‌های موردانتظار (تعداد زنجیره‌ها در پیکره‌آزمون) یکسان هستند، دقت و فراخوانی یکی بوده و عملکرد برنامه با صحت اندازه‌گیری می‌شود:

$$Accuracy = \frac{\text{Number of correct outputs proposed by program}}{\text{Total number of English queries in the test set}} = \frac{75}{100} = 75\%$$

نتایج به دست آمده از معادله‌ی خودکار با همان مجموعه از واژه‌ها که از جفت جملات انگلیسی و فارسی در پیکره به صورت دستی در سطح واژه هم‌تراز شده بودند مقایسه شد. نتایج آزمایش نشان داد که برنامه مطابق ما صحتی معادل ۷۵ درصد بدست آورده که بسیار امیدوارکننده است. طبیعی است در حالتی که این الگوریتم بر روی یک پیکره طبیعی از زبان انگلیسی که در آن انواع مختلف واژه‌ها (هم واژه‌های محتوایی و هم واژه‌های دستوری) یافت می‌شوند اعمال شود، صحت این روش مسلماً تا حد زیادی افزایش خواهد یافت. شکل (۳) نمونه‌ای از یک مطابقت موازی در سطح جمله که شامل ۵ جفت جمله با جستار انگلیسی *campaign*، معادل‌های فارسی آن و مقدار اطلاعات متقابل آن‌ها (در ستون چهارم) است را نشان می‌دهد

No	English Sentence	Farsi Sentence	MI value
1	double standards destroy the moral superiority of actually-existing democracies and allow despotic rulers everywhere in the world, including the middle east and iran, to present a better image of themselves in their propaganda campaign against democratic forces.	معیارهای دوگانه برتری اخلاقی دموکراسی‌های واقعاً موجود را از میان می‌برد، و به حاکمان خودکامه در سراسر جهان از جمله منطقه‌ی خاورمیانه و نیز ایران این امکان را می‌دهد که بتوانند در مبارزه تبلیغاتی علیه نیروهای دموکرات خود را در وضوح و بهتری قرار دهند.	7.967
2	the united states pressed russia on friday to join the west in rebuking iran at the united nations as part of a u.s. - led campaign to curb the islamic republic's nuclear program .c	امریکا روز جمعه روسیه را زیر فشار گذاشت تا در برخورد با ایران در سازمان ملل، به عنوان پشتیبان مبارزه امریکا برای متوقف سازیدن برنامه هسته ای جمهوری اسلامی غرب، همسو گردد.	7.967
3	he stipulated that iran is decisive in its campaign against terrorism and in this regard takes action on the basis of its international responsibilities, and in the future also acts on the same basis.	دکتر آصفی تصریح کرد: ایران در مبارزه با تروریسم کاملاً فاطم بوده و بر اساس مسئولیتهای بین المللی خود اقدام کرده و در آینده نیز بر همین اساس عمل خواهد کرد.	7.967
4	dr assefi said the islamic republic of iran emphasizes on the non-selective, explicit and transparent campaign against the inhuman phenomenon of terrorism, and we believe that eradication of terrorism will be realized only by participation of members of international community and through a collective action	دکتر حمیدرضا آصفی اضافه کرد: مبارزه غیر گزینشی، صریح و شفاف با پدیده همد انسانیت تروریسم مورد تأکید ایران می باشد و ما معتقدیم ریشه کنی تروریسم تنها با مشارکت اعضای جامعه جهانی و در نتیجه یک اقدام جمعی مؤثر واقع خواهد شد.	7.967

شکل ۳) نمونه‌ای از یک مطابقه موازی در سطح جمله

۶. نتیجه‌گیری و چشم‌انداز آینده

بی‌شک مطابقه‌های دوزبانه منابع ارزشمندی برای زبانشناسان، مترجمان و نیز کاربران رایانه در حوزه پردازش متن و بویژه بازیابی اطلاعات هستند. این مطابقه‌ها از مجموعه حجیمی از جملات همتراز شده در سطح جمله به‌عنوان پایگاه داده‌ای قابل جستجو استفاده می‌کنند. کاربر می‌تواند یک واحد زبانشناختی یا زنجیره‌ای از واحدها در یک زبان را جستجو کند و مطابقه می‌تواند تمام جملات به آن زبان که در آن‌ها واحد (های) مورد جستجو ظاهر شده‌اند همراه با جملات متناظر آن‌ها به زبان دیگر را نمایش می‌دهد.

در این مقاله تلاش بر این بود تا روشی نسبتاً نوین برای همترازسازی خودکار مجموعه‌ای از جملات موازی با استفاده از رویکردی آماری بنام آماره اطلاعات متقابل ارائه شود. برای این منظور یک پیکره موازی انگلیسی فارسی که به‌طور دستی در سطح جمله همترازسازی شده است مورد استفاده قرار گرفت. مقدار اطلاعات متقابل نشان‌دهنده درجه همبستگی معنایی بین واژه‌ها است.

همترازسازی در سطح واژه مطابقه‌های دوزبانه را غنی‌تر می‌سازد و کاربردهای گوناگونی از جمله ترجمه ماشینی آماری، بازیابی اطلاعات دوزبانه، یادگیری زبان، واژه‌نگاری رایانه‌ای و مانند آن دارد. روش ارائه‌شده در این تحقیق روشی مستقل از زبان است که صرفاً تکیه بر مقیاس‌های آماری دارد. از این‌رو، الگوریتم این روش می‌تواند برای هر جفت زبان دیگری که پیکره موازی آن‌ها موجود باشد بکار رود.

هنوز هم اصلاحات زیاد دیگری به این سیستم اضافه شود تا به تولید برون‌دادهای بهتری منجر شود. یکی این که نویسندگان قصد دارند پایگاه داده‌ای پیکره موجود را با اطلاعات جدید بیشتری بطور مداوم بروز رسانی نمایند بطوری که مطابقه موازی موجود که در سطح جمله قادر به معادل‌یابی است به مطابقه قدرتمندتری تبدیل شود. به‌علاوه، تفکیک حوزه‌های موضوعی نیز می‌تواند به بازدهی بیشتر و مؤثرتری منجر شود.

قردانی

این پژوهش با استفاده از اعتبارات دانشگاه پیام نور در قالب طرح پژوهشی انجام شده است.

منابع

- نظارات، امین و موسوی میانگه، طیبه (۱۳۹۰). طراحی و پیاده‌سازی یک سامانه بازیابی اطلاعات دوزبانه با استفاده از پیکره‌های زبانی. *پردازش و مدیریت اطلاعات، ویژه نامه ذخیره، بازیابی و مدیریت اطلاعات: ۲۱۲-۱۹۷*.
- Barlow, M. (2002). ParaConc: Concordance software for multilingual parallel corpora. In: *Proceedings of the Third International Conference on Language Resources and Evaluation. Workshop on Language Resources in Translation Work and Research.*, Las Palmas, Spain, pp. 20-24.
- Brown, P. F., Della Pietra, V. S. A.; Della Pietra, V. J.; and Mercer, R. L. (1993). The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics* 19 (2), 263° 312.
- Church, K. W. and Hanks, P. (1990). Word association norms, mutual information, and lexicography. *Computational Linguistics* 16 (1) , 22-29.
- Fung, P. (1995). A pattern matching method for finding noun and proper noun translations from noisy parallel corpora. *Proc. of the 33rd Annual Meeting of the ACL*, 236-243.
- Gale, .. A. and K. .. Church, (1991). Identifying word correspondences in parallel texts *Proceedings of the 4th DARPA Speech and Natural Language Workshop::152-157*, Pacific Grove, CA.
- Dagan, I. Kenneth W. Ch, and William A. G. (1993). Robust bilingual word alignment for machine-aided translation. In: *Proc. of the Workshop on Very Large Corpora: Academic and Industrial Perspectives*.
- Inoue, N. and Nogaito, I. (1993). Automatic construction of the Japanese-English dictionary from bilingual text. *Technical Report of IEICE, NLC: 39-93*.
- Ishimoto, H and Nagao, M. (1994). Automatic construction of a bilingual dictionary of technical terms from parallel texts. *Technical Report of IPSJ, NL: 102-11*.
- Johns, T. (1986). Microconcord: A language-learner's research tool. *System* 14(2), 151-162.
- Johns, T. (1998). Multiconcord: the lingua multilingual parallel concordancer for windows. Available on: http://web.bham.ac.uk/johnstf/l_text.htm. Accessed Feb 03
- Kaji, H and Aizoni, T. (1996). Extracting word correspondences from bilingual corpora based on word co-occurrences information. *The 16th International Conference on Computational Linguistics*, pp. 23 ° 28. Copenhagen, Denmark
- Kumano, A. and Hirakawa, H. (1994). Building an MT dictionary from parallel texts based on linguistic and statistical information. *Proc. of COLING'94: 76-81*.
- Kupiec, J. (1993). An algorithm for finding noun phrase correspondences in bilingual corpora. *Proceedings of the 31st Annual Meeting of the ACL: 17-22*, Columbus, Ohio.
- Mausser, A., Matusov, E and Ney, H. (2006). Training a Statistical Machine Translation System without GIZA++. *International Conference on Language Resources and Evaluation (LREC): 715-720*, Genoa, Italy.
- Mosavi Miangah, T. (2006). Applications of corpora in translation. *Translation Studies* 12, 43-56.
- Mosavi Miangah, T. (2009). Constructing a large-scale English-Persian Parallel Corpus. *META* 54 (1), 181-188.
- Mosavi Miangah (in Press). FarsiTag: A part of speech tagging system for Persian. *Journal of Quantitative Linguistics*.
- Och, Franz J. (2000). Giza++: Training of statistical translation models. Available at: <http://www-i6.informatik.rwthachen.de/~och/software/GIZA++.html>.
- Och, F. J. and Ney, H. (2003). A systematic comparison of various statistical alignment models. *Computational Linguistics* 29(1), 1-19.
- Och, F. and Ney, H. (2001). Improved Statistical Alignment Models, *Proceedings of ACL 2001*.

- Rensik, .. (1998). PParallel strands: A preliminary investigation into mining the web for bilingual text.. In *Proceedings of the Third Conference of the Association for Machine Translation in the Americas, AMTA-98, in Lecture Notes in Artificial Intelligence, 1529*: 28-31, Langhorne, PA, October,.
- Rensik, .. (1999): Mining the web for bilingual text . *Proc. of 37th Meeting of the ACL*. Maryland: 527-534.
- Scott, M. (2000). *WordSmith Tools Version 3.0* [Computer software]. Oxford: Oxford University Press.
- Simard, M. and Plamondon, P. (1998). Bilingual sentence alignment: balancing robustness and accuracy. *Machine Translation* 13, 59° 80.
- Wang L. (2001). Exploring parallel concordancing in English and Chinese. *Language Learning & Technology* 5(3), 174-184.
- Wu, J. C.; Yeh, K.; Chuang, C.; Thomas C., Shei, W. C. and Chang, J. (2003). TotalRecall: A bilingual concordance for computer assisted translation and language learning. *Association for Computational Linguistics*: 201-204.
- Yamamoto, Y. and Sakamoto, M. (1993). Extraction of technical term bilingual dictionary from bilingual corpus. *Technical Report of IPSJ, NL*: 12-94. (in Japanese).



