

دانشگاه فردوسی مشهد
دانشکده مهندسی - گروه مهندسی کامپیوتر

پایان نامه کارشناسی ارشد

ارائه رهیافتی جدید برای تولید پیکره موازی انگلیسی-فارسی

تهیه و تنظیم:

سید احمد جکیان طوسی

استاد راهنما:

دکتر محسن کاهانی

استاد مشاور:

دکتر هادی صدوقی یزدی

تابستان 91

تقدیم بہ:

بہ ہمہ معانم، کہ سرچشمہ لطیف علم اندوپیامبران نش فرہنگی ام

بہ پدرم، کہ استاد معرفت است واسوہ شکیبانی ام

بہ مادرم کہ معدن محبت است و سرایہ آسمانی ام

بہ ہمسرم کہ مرواریدی بدیل عشق است و امید بخش زندگانی ام

و بہ فرزندان عزیز و عزیزان دلہام

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

تعهدنامه

اینجانب سید احمد جکیان طوسی دانشجوی کارشناسی ارشد رشته مهندسی کامپیوتر، گرایش هوش مصنوعی، دانشکده مهندسی، دانشگاه فردوسی مشهد، نویسنده پایان‌نامه "ارائه رهیافتی جدید برای تولید پیکره های موازی انگلیسی فارسی" تحت راهنمایی دکتر محسن کاهانی متعهد می‌شوم:

تحقیقات در این پایان‌نامه توسط اینجانب انجام شده و از صحت و اصالت برخوردار است.

در استفاده از نتایج پژوهش‌های محققان دیگر به مرجع مورد استفاده استناد شده است.

مطالب مندرج در پایان‌نامه تاکنون توسط خود و یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است.

کلیه حقوق معنوی این اثر متعلق به دانشگاه فردوسی مشهد می‌باشد و مقالات مستخرج با نام "دانشگاه فردوسی مشهد" و یا "Ferdowsi University of Mashhad" به چاپ خواهد رسید.

حقوق معنوی تمام افرادی که در به دست آمدن نتایج اصلی پایان‌نامه تأثیرگذار بوده‌اند در مقالات مستخرج از رساله رعایت شده است.

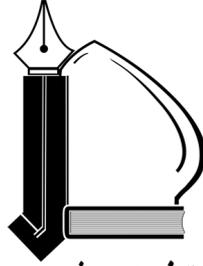
در کلیه مراحل انجام این پایان‌نامه، در مواردی که از موجود زنده (یا بافت‌های آن‌ها) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است. در کلیه مراحل انجام این پایان‌نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است، اصل رازداری، ضوابط و اصول اخلاق انسانی رعایت شده است.

امضای دانشجو

تاریخ

مالکیت نتایج و حق نشر

کلیه حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه‌های رایانه‌ای، نرم‌افزارها و تجهیزات ساخته شده) متعلق به دانشگاه فردوسی مشهد می‌باشد. این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود. استفاده از اطلاعات و نتایج موجود در پایان‌نامه بدون ذکر مرجع مجاز نمی‌باشد.



دانشگاه فردوسی مشهد
دانشکده مهندسی - گروه مهندسی کامپیوتر

آزمایشگاه تخصصی فناوری وب WTLAB

پایان نامه کارشناسی ارشد

ارائه رهیافتی جدید برای تولید پیکره موازی انگلیسی-فارسی

تهیه و تنظیم:

سید احمد جکیان طوسی

استاد راهنما:

دکتر محسن کاهانی

استاد مشاور:

دکتر هادی صدوقی یزدی

تابستان 91

تقدیر و تشکر

پس از حمد بیکران از لطف و رحمت الهی و پیش از ارایه هر مطلبی، لازم می‌دانم با احترام از کسانی نام ببرم که در طول انجام پژوهش‌های مربوط به پایان نامه دوره کارشناسی ارشد، مرایاری نمودند و با کمک رسانی مخلصانه خویش امکان غلبه بر موانع و چالش‌های موجود را فراهم کردند.

میش از همه، شایسته می‌دانم، تشکر ویژه خود را تقدیم به استاد فرهیخته و فرزانه گران قدر، دکتر محسن کاغذی‌نایم که در این مدت، با گذر خویش رادلسوزانه و صبورانه حمایت نموده و با توفیق‌های مؤثر خود، عامل مهمی در پیشرفت امور علمی ام بوده‌اند.

پس از استاد مشاور پرورده، دکتر هادی صدوقی یزدی که سهم بسزایی در طی نمودن صحیح این مسیر داشتند، به طور ویژه، تقدیر و تشکر می‌نمایم.

از استاد فرزانه زبان‌شناسی، پروفسور نادجما نگیری نیز به دلیل همراهی موثرشان کمال سپاس و قدردانی را دارم.

بدون تردید اگر همکاری سایر دوستان و همکارانم در آزمایشگاه فناوری وب نبود، انجام کارهای مربوط به جمع‌آوری داده و پیاده‌سازی سیستم با مشکلات بسیار همراه می‌شد به همین منظور با تقدیم سپاس به مجموعه توانمند آزمایشگاه فناوری وب شایسته می‌دانم به طور ویژه از زحمات دوستان عزیزم آقایان مهندس، رضا سعیدی، نیاصالحی، احمد استیری، امین العیدیار و بهداد مهدی به جهت همکاری‌هایی که در مقطع کوناگون با اینجانب داشته‌اند، صمیمانه تشکر می‌نمایم.

در نهایت با امید به بهبود هر چه بیشتر فرآیندهای تحقیقاتی، توفیقات روز افزون را برای همه فعالان عرصه رشد و شکوفایی علمی، در اسکله مهندسی دانشگاه فردوسی مشهد آرزو می‌کنم.

سید احمد طوسی

چکیده:

در این پژوهش، برای اولین بار مدلی ترکیبی برای تراز بندی جملات، جهت ساخت پیکره‌های موازی انگلیسی-فارسی ارائه شده است. در حالت کلی چارچوب روش پیشنهادی، غیر وابسته به زبان‌های مبدأ و مقصد بوده و از آن می‌توان برای تولید پیکره‌های موازی، برای هر جفت زبان دیگر، نیز استفاده کرد. نتایج بدست آمده از پیاده‌سازی‌ها نشان می‌دهد که بکار بردن ویژگی‌های زبانی و غیر زبانی، عملکرد سیستم را تا حد قابل قبولی بهبود خواهد بخشید. همچنین در این تحقیق میزان تأثیر استفاده از شباهت‌های طولی، ترجمه تحت‌اللفظی و نقش دستوری کلمات، به صورت مستقل یا ترکیبی، برای عملیات تراز بندی مورد بررسی قرار گرفته است. بکارگیری طبقه بندهای سلسله‌مراتبی در تشخیص نوع تراز بندی به عنوان یکی از شاخص‌های اصلی سیستم پیشنهادی محسوب می‌شود که موجب بالا رفتن دقت و سرعت عملیات تولید پیکره نسبت به سایر مدل‌ها می‌شود. از ویژگی‌های دیگر این روش، قابلیت توسعه پذیری آن است. به این ترتیب که می‌توان با گنجاندن خصوصیات (در بدنه بردارهای ورودی) که ممکن است در آینده برای تشخیص بهتر نوع تراز بندی مورد توجه قرار گیرند کیفیت سیستم تراز بندی را ارتقاء بخشید. با این حال، چالش اساسی این رهیافت و بسیاری از روش‌های پیشین، وجود سلايق متنوع در ترجمه متون است. این مسئله در برخی موارد موجب تولید جملاتی می‌گردد که مشابهت آن‌ها با متن اصلی در حد معنا بوده و تنها قابل درک و تشخیص برای انسان می‌باشد. در چنین وضعیتی کار استخراج جفت عبارات معادل، بسیار دشوار است. علاوه بر این در مورد متونی که مقید به رعایت ساختار دستوری زبان فارسی و انگلیسی نمی‌باشند عملیات تشخیص جملات هم تراز ممکن است به خوبی انجام نپذیرد.

کلمات کلیدی:

تراز بندی جملات، پیکره موازی، ترجمه آماری

فهرست مطالب

1	1- مقدمه
2	1-1- انگیزه
4	2-1- راه حل پیشنهادی
6	3-1- ساختار پایان نامه
9	2- مرور ادبیات
10	1-2- ترجمه ماشینی
12	2-2- ترجمه آماری
14	3-2- پیکره دو زبانه
14	4-2- پیکره موازی
16	5-2- پیکره تطبیقی
16	6-2- تراز بندی جملات
17	7-2- پیش پردازش زبانی
18	8-2- برچسب زنی دستوری کلمات
18	9-2- ریشه یابی کلمات
19	10-2- طبقه بندی
20	11-2- طبقه بند تک کلاسه
20	12-2- طبقه بند دو کلاسه
21	13-2- طبقه بند چند کلاسه
23	14-2- آموزش و اعتبار سنجی متقاطع n بخشی
24	15-2- درخت تصمیم

25 CART	16-2
26 SVM	17-2
28 C-SVM	18-2
30	19-2
31	1-19-2
32	2-19-2
33 F	3-19-2
33	20-2
49	21-2
50	1-21-2
50	2-21-2
50	3-21-2
50	4-21-2
50 TEP	5-21-2
51	22-2
58	23-2
60	3-3
60	1-3
61	2-3
63	3-3
67	4-3

71	مدل L	1-4-3
72	مدل P	2-4-3
73	مدل T	3-4-3
76	مدل‌های ترکیبی LP,PT,LT,LPT	4-4-3
77	طبقه بند سلسله‌مراتبی	5-3
80	الگوریتم کلی تراژبندی	6-3
82	خلاصه فصل	7-3
84	پیاده سازی و ارزیابی	4
84	مقدمه	1-4
86	پیاده سازی	2-4
87	گرد آوری متون موازی	1-2-4
89	تولید پیکره آموزشی FEP	2-2-4
92	حذف حروف و علائم زائد	3-2-4
94	برچسب‌زنی دستوری پیکره	4-2-4
95	نگاشت برچسب‌های دستوری	5-2-4
96	اصلاح‌گر برچسب‌های دستوری	6-2-4
98	لغتنامه انگلیسی فارسی	7-2-4
98	تعلیم و آزمون سیستم	8-2-4
99	ارزیابی سیستم	3-4
100	بررسی تأثیر استفاده از پیکره FEP و TEP	1-3-4
101	بررسی تأثیر به‌کارگیری انواع بردار ویژگی	4-4

102	5-4- بررسی نقش عملکرد انواع طبقه بندها در کارایی سیستم.....
106	6-4- بررسی تأثیر اصلاح گر برچسب‌های دستوری.....
107	7-4- بررسی اصلاح مدل مؤلفه های مدل T در کارایی سیستم.....
108	8-4- ارزیابی مقایسه ای سیستم پیشنهادی.....
111	9-4- ارزیابی تأثیر اصلاحات ساختاری در سیستم پیشنهادی.....
115	10-4- خلاصه.....
118	5- نتیجه گیری و کارهای آتی.....
118	5-1- نتیجه گیری.....
119	5-2- کارهای آینده.....
123	6- مراجع.....

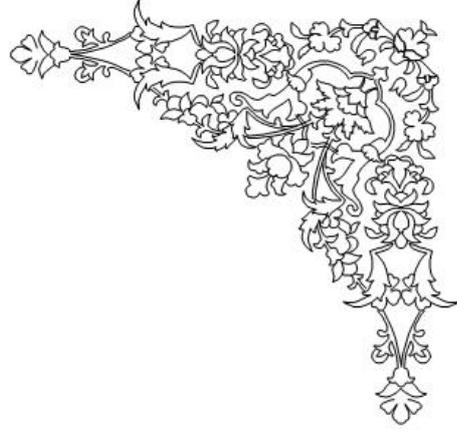
فهرست نمودارها و شکل‌ها

- شکل (1-2) نحوه کار سیستم مترجم آماری 14
- شکل (2-2) مراحل انجام اعتبارسنجی متقاطع 4 بخشی 24
- شکل (3-2) نمایی از یک درخت تصمیم 24
- شکل (4-2) نمایش نموداری نتایج واقعاً درست و نتایج درست تشخیص داده شده 32
- شکل (5-2) نمونه ای از پیکره دوزبانه پارلمان کانادا 34
- شکل (6-2) نسبت طول پاراگراف‌های موازی انگلیسی - آلمانی 35
- شکل (7-2) نمودار همبستگی طول پاراگراف‌های انگلیسی با مجذور اختلاف طول پاراگراف معادل در زبان آلمانی ... 37
- شکل (8-2) نحوی شکستن متن انگلیسی و فارسی براساس موقعیت شروع جملات 55
- شکل (9-2) نمودار توزیع طول جملات انگلیسی بر حسب تعداد حروف 57
- شکل (10-2) نمودار توزیع طول جملات فارسی بر حسب تعداد حروف 57
- شکل (1-3) ساختار کلی مدل پیشنهادی - مرحله تعلیم سیستم 62
- شکل (2-3) ساختار کلی مدل پیشنهادی - مرحله تولید پیکره موازی 63
- شکل (3-3) نمودار کلی محاسبه شباهت دو عبارت فارسی و انگلیسی از لحاظ ترجمه‌ای (مدل T) 74
- شکل (4-3) انواع بردارهای ویژگی 76
- شکل (5-3) ساختار کلی طبقه‌بند سلسله‌مراتبی مدل پیشنهادی 77
- شکل (6-3) ساختار کلی طبقه‌بند چندکلاسه MDT 78
- شکل (1-4) نمودار فرآیندهای انجام شده در جریان اجرای پروژه 85
- شکل (2-4) نمایی از صفحه اصلی نرم افزار ایجاد پیکره 91
- شکل (3-4) نمایی از صفحه اصلی نرم افزار الحاق‌گر 91
- شکل (4-4) نمایی از ابزار توسعه یافته برچسب‌زن دستوری فارسی در گیت 95
- شکل (5-4) اصلاح‌گر برچسب‌های دستوری به روش نیمه خودکار 97
- شکل (6-4) نمودار ارزیابی تأثیر کیفیت پیکره آموزشی بر روی سیستم ترازبند MDT 101
- شکل (7-4) نمودار مقایسه عملکرد انواع طبقه‌بند ها در سیستم های پیشنهادی 104
- شکل (8-4) نمودار تأثیر اصلاح برچسب‌های دستوری در مدل P 107

- شکل (4-9) نمودار تأثیر استفاده از مترجم آنلاین گوگل بر کارایی مؤلفه های مدل T 108
- شکل (4-10) مقایسه کیفیت تشخیص هم ترازى برنامه haling و نسخه پياده سازى شده LPT بر روى پيكره FEP 109
- شکل (4-11) ارزیابی کیفیت روش گیل در میان انواع مدل‌های پروژه پیشنهادی 110
- شکل (4-12) نمودار مقایسه عملکرد سیستم های پیشنهادی 112
- شکل (4-13) ساخت نسخه های مختلف پیکره آموزشی FEP 115

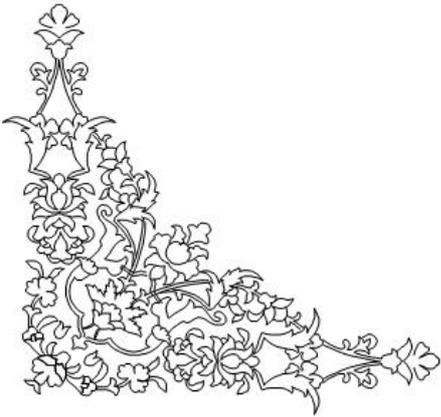
فهرست جدول‌ها

- جدول (1-2) مقادیر مربوط به محاسبه احتمالات پیشین [GAL93] 37
- جدول (2-2) اندازه گیری خطای روش گیل بر روی پیکره انگلیسی فرانسوی و انگلیسی آلمانی 40
- جدول (3-2) مقادیر p(a) بر اساس مستندات پارلمان کانادا [GAL93] 43
- جدول (4-2) نتایج روش ترکیبی مبتنی بر طول و مبتنی بر هم‌ریشه‌ها [GAL93] 44
- جدول (5-2) مقایسه عملکرد GMM در برابر روش‌های مبتنی بر طول و روش ترکیبی 48
- جدول (6-2) نمونه‌هایی از واژگان محاوره ای موجود در پیکره TEP 51
- جدول (7-2) نتایج روش استخراج جملات هم‌تراز از درون مخزن ویکی‌پدیا [MOH10] 54
- جدول (8-2) نمونه ای از زیر نویس های فارسی و انگلیسی 57
- جدول (1-3) نمونه ای از پیکره آموزشی تراز شده به همراه برچسب هم‌ترازی 69
- جدول (2-3) مجموعه جملات انتخاب شده جهت ساخت اولین بردار ویژگی 70
- جدول (1-4) نسبت توزیع مطالب جمع آوری شده جهت استفاده در پیکره آموزشی 89
- جدول (2-4) حجم پیکره FEP به تفکیک موضوعات 92
- جدول (3-4) پاک‌سازی پیکره آموزشی از حروف زائد 93
- جدول (4-4) جداکننده جملات فارسی و انگلیسی 94
- جدول (5-4) نسبت وجود انواع هم‌ترازی ها در پیکره آموزشی 94
- جدول (6-4) مقایسه نسبت به کارگیری انواع برچسب‌های دستوری در هر جمله فارسی و انگلیسی ... 97
- جدول (7-4) نتایج آزمون سیستم به روش اعتبار سنجی متقاطع 10 بخشی در مدل MDT 99
- جدول (8-4) مقایسه نتایج سیستم ترازبندی با دو پیکره FEP , TEP براساس معیار اندازه F 100
- جدول (9-4) مقایسه تأثیر انواع بردار ویژگی بر روی شاخص‌های کیفی سیستم 101
- جدول (10-4) نتایج توزیع نمونه های پیکره توسط الگوریتم‌های خوشه‌بندی 103
- جدول (11-4) مقایسه عملکرد انواع طبقه بند 105
- جدول (12-4) مقایسه عملکرد سیستم های پیشنهادی 112
- جدول (13-4) جدول مقایسه زمان اجرای سیستم پیشنهادی و C-SVM 113
- جدول (14-4) حجم و زمان تولید بردارهای ویژگی 114



فصل اول:

مقدمه



1- مقدمه

بنابر آمار منتشر شده از سوی یکی از نشریات معتبر مربوط به زبان‌شناسی دنیا [VER04]، ترجمه ماشینی به عنوان یک نامزد جدی در تصاحب اولین جایگاه در میان 10 فناوری برتر دنیا مطرح می‌باشد. این فناوری‌های نوین مدعی‌اند که می‌توانند دنیا را به طور اساسی دچار تحول نمایند. در این فهرست از فناوری‌های نانو، بیوتکنولوژی، تکنولوژی اطلاعات، علوم شناختی، رباتیک و هوش مصنوعی نیز به عنوان دیگر فناوری‌های برتر نام برده شده است.

تاکنون دانشمندان و نظریه‌پردازان علوم مربوط به زبان‌شناسی و فنون مهندسی رایانه، راهکارهای مشترک و متنوعی را برای حل مسئله ترجمه ارایه نموده‌اند. با ظهور پدیده اینترنت و برداشته شدن موانع ارتباطی میان ملت‌ها، این موضوع بیش از پیش، مورد تقاضای ملل گوناگون، که دارای گویش‌ها و فرهنگ‌های مختلف هستند، قرار گرفته است. با اینکه بین‌المللی کردن زبان انگلیسی، به عنوان یک راهکار اولیه، برای تسهیل ارتباطات انسان‌ها در سراسر دنیا، مورد قبول جامعه جهانی است، اما وفاداری ملت‌ها به حفظ فرهنگ گفتاری و نوشتاری خود مانع از رفع کامل نیازشان به سیستم‌های مترجم می‌باشد. از این رو از دیرباز، با وجود دوران رخوت و سستی در طراحی این سیستم‌ها باز هم این موضوع فراموش نشده و در مقاطع گوناگون به عنوان یکی از معمای علمی ممتاز مورد توجه محققان بوده است.

1-1- انگیزه

یک مترجم ایده‌آل، در حقیقت یک انتقال دهنده دانش است که این کار را با تبدیل معنایی عبارات مبدأ به جملاتی در زبان مقصد انجام می‌دهد. به بیان ساده تر یک مترجم خوب مترجمی است که بتواند منظور نویسنده یا گوینده را که در قالب زبان خاص بیان شده است به صورت جملاتی به زبان خاص دیگر که متعلق به خواننده یا شنونده است تبدیل نماید. با وجود سادگی این تعریف، پیاده‌سازی چنین سیستمی نه تنها ساده نیست بلکه بسیار پیچیده و نیازمند دقت در ابعاد مختلف می‌باشد. از این رو، انگیزه بسیاری از محققان فنون ترجمه و علوم زبان شناسی رایانه ای یافتن راهی است که از آن طریق بتوان به ترجمه ایده آل برای متون مختلف دست پیدا کرد.

یکی از راه‌های ساخت مترجم ماشینی، استفاده از دانش آماری نهفته در متون ترجمه شده‌ای است که داده های آن در قالب پیکره های موازی برای زبان مبدأ و مقصد در اختیار می‌باشند. [RAM04]. به عبارت دیگر، ایده اصلی این روش - که با عنوان ترجمه ماشینی آماری شناخته می‌شود- انجام عملیات ترجمه بر اساس یافتن شبیه‌ترین عبارات (موجود در) پیکره موازی به عبارت ورودی می‌باشد. صرف نظر از تکنیکی که برای پیدا کردن زیر عبارات مشابه مورد استفاده قرار می‌گیرد به نظر می‌رسد مهم‌ترین چالش موجود در این تکنیک، دسترسی به پیکره موازی به حد کافی بزرگ برای دو زبان مبدأ و مقصد می‌باشد.

تا کنون روش‌های مختلفی برای جمع آوری متون موازی از روی نوشته های دو زبانه پیشنهاد شده است. ولی به این دلیل که همواره مترجم انسانی به دلایلی ترجمه خود را جمله به جمله انجام نمی‌دهد، مسئله دیگری با عنوان تراز بندی مطرح می‌شود. به این معنا که ممکن است مثلاً مترجم یک جمله زبان مبدأ را به دو یا چند جمله در زبان مقصد و یا بالعکس ترجمه نماید. از این رو در سال‌های اخیر، مقالات متعددی با عنوان تراز بندی جملات به منظور تولید پیکره های موازی ارایه شده‌اند که با دیدگاه‌های متفاوت به دنبال تشخیص نوع تراز بندی بکار رفته در ترجمه می‌باشند.

به عنوان نمونه، در برخی از روش‌های تراز بندی جملات، بدون توجه به ویژگی‌های زبانی و معنایی دو زبان، تنها بر اساس تعداد کاراکترهای تشکیل دهنده جمله/جمله‌ها، عمل تراز بندی انجام می‌پذیرد. ایده اساسی

این رویکرد را می‌توان در یک عبارت چنین خلاصه نمود که معمولاً عبارات با طول کم به عبارات با طول کم و عبارات با طول بیشتر به عبارات با طول بیشتر ترجمه می‌شوند. [GAL93]

بررسی‌ها نشان می‌دهد که در مورد برخی از جفت زبان‌ها، نقش دستوری کلمات در هنگام ترجمه چندان دست‌خوش تغییر نمی‌شوند. یعنی واژه‌هایی که با نقش فعل در جمله مبدأ بکار رفته‌اند در عبارت ترجمه نیز به صورت فعل ظاهر می‌شوند. از این رو ابزارهایی برای تولید برچسب‌های زبانی چه از نوع معنایی و چه از نوع دستوری مورد توجه بخشی دیگری از محققان قرار گرفت. [CHE95]

در تعداد دیگری از روش‌های مطرح شده، [MEY98] از فرهنگ واژگان برای انجام عمل تراز بندی بهره گرفته شده است. در واقع لغت نامه نقش یک مترجم واژه به واژه را در این روش ایفا می‌کند که کار یافتن شبیه‌ترین عبارت/عبارات را در دو متن موازی انجام می‌دهد.

عده دیگری از پژوهشگران، با توجه به رشد روزافزون وب و انتشار سریع داده‌های متنی در آن، تکنیک‌هایی را برای یافتن جفت عبارات موازی از روی وب مطرح نموده‌اند. به عنوان مثال در سایت‌های خبری بین‌المللی یا دایره‌المعارف‌های چند زبانی، معمولاً ابر پیوندها¹ در دو عبارت موازی، به یک محل مشترک اشاره می‌کنند. تا کنون بوسیله این روش ابتکاری، تعداد قابل توجهی عبارات موازی، از درون مخازن وب استخراج شده است. [MOH10]

در این میان برخی روش‌های ترکیبی [ASW05] مطرح شده‌اند که انگیزه اصلی آن‌ها از ترکیب روش‌های گوناگون، پوشش نقاط ضعف برخی از روش‌ها، به وسیله روش‌های دیگر می‌باشد. نتایج بررسی‌ها نشان دهنده این موضوع است که روش‌های ترکیبی نسبتاً بهتر از روش‌های قبلی عمل می‌کنند.

تا کنون عملیات پیاده سازی روش‌هایی که پیش از این به آن‌ها اشاره اجمالی شد بر روی تعداد بسیاری از جفت زبان‌های دنیا صورت گرفته است که نتیجه آن تولید پیکره حجیم موازی بوده است. همانطور که گفته شد استفاده از این پیکره‌های عظیم در مدل‌های ترجمه آماری سبب ارتقاء کیفیت ترجمه‌های ماشینی بدست آمده شده است. با این وجود، متأسفانه سهم پارسی‌زبانان در تولید خودکار پیکره انگلیسی

¹ Hyperlink

فارسی چندان قابل توجه نبوده است. البته در حال حاضر حرکت‌های مثبت توسط تیم‌های دانشگاهی مرتبط در حال انجام می‌باشد لکن اهمیت موضوع به اندازه‌ای است که مستعد انجام کوشش‌های بیشتر بوده و رسیدگی به این امر را به یک ضرورت قطعی تبدیل نموده است.

انگیزه اصلی پایان نامه حاضر ارایه راهکاری مؤثر، برای تولید خودکار پیکره موازی انگلیسی فارسی می‌باشد. رهیافت کلی تحقیق پیش رو، مبتنی بر یک روش ترکیبی، می‌باشد. چارچوب کلی روش ترازبندی پیشنهاد شده، مستقل از هر دو زبان مبدأ و مقصد می‌باشد. با این حال امکان توجه به خصوصیات زبانی و فرا زبانی به صورت مستقل و یا با هم نیز وجود دارد. نتایج حاصل از کاربرد روش پیشنهادی برای تولید پیکره‌های حجیم، می‌تواند باعث توسعه سیستم‌های مترجم قوی‌تر شود و لذا تأثیر بسزایی در تسهیل فرآیندهای تحقیقاتی، آموزشی و تجاری خواهد داشت.

1-2- راه حل پیشنهادی

تا کنون روش‌های متنوعی برای تولید پیکره های حجیم موازی ارائه شده است. تعدادی از روش‌های مذکور از معیارهای فرا زبانی^۲ یعنی معیارهای غیر وابسته به زبان خاص، برای حل مسئله تراز بندی عبارات استفاده می‌کنند. این ویژگی سبب شده است که از آن بتوان برای تولید هر نوع پیکره دو زبانی استفاده نمود.

در گروه دیگری از روش‌ها به ویژگی‌های زبان‌های مطرح در پیکره، توجه خاص شده است. از این رو نمی‌توان انتظار داشت که این گونه مدل‌ها برای هر جفت زبان موجود در دنیا کاربرد داشته باشند. مزیت کلی این گروه از روش‌ها را می‌توان به دلیل بالا بودن دقت آن‌ها در تشخیص انواع تراز بندی جملات در دو متن موازی دانست.

در مدل پیشنهادی که در این تحقیق به آن پرداخته شده است هم به خصوصیات فرا زبانی و هم به ویژگی‌های زبانی توجه شده است، تا ضمن برخورداری از مزایای روش‌های فرا زبانی، با بهره‌مندی از ویژگی‌های خاص زبانی، امکان تولید پیکره های مطلوب انگلیسی فارسی به وجود آید. بر این اساس

² Language-independent

می‌توان از این مدل با اطمینان برای ترازبندی جملات مربوط به هر یک از زبان‌های هم خانواده با زبان فارسی از قبیل هندی و عربی در برابر عبارات انگلیسی استفاده نمود.

همان‌طور که در ادامه نیز تشریح خواهد شد استفاده از تئوری‌های یادگیری ماشین از قبیل طبقه‌بندهای درخت تصمیم^۳ CART و ماشین بردار پشتیبان^۴ در تراز بندی جفت جملات انگلیسی و فارسی جزء مهم‌ترین بخش‌های سیستم پیشنهادی می‌باشد. علاوه بر این‌ها استفاده از ابزار زبان شناسی و انجام برخی اصلاحات در آن‌ها به منظور کاربرد بهتر در ساختار مدل ترازبندی، از دیگر امتیازات این روش می‌باشد.

در یک نگاه کلی، روشی که در این پایان نامه به آن اشاره می‌شود به این صورت است که ابتدا یک پیکره تراز شده با حجم و دقت قابل قبول به یک سیستم یادگیرنده آموزش داده می‌شود. پس از آن، سیستم توانایی انجام عملیات ترازبندی، بر روی عبارات جفت متن موازی را خواهد داشت.

مهم‌ترین ویژگی مدل پیشنهاد شده در این پایان نامه، عدم وابستگی آن به جفت زبان‌های مربوط به پیکره - در حالت کلی - می‌باشد. البته در صورت در دسترس بودن هر یک از خصوصیات زبانی موثر در ترازبندی، می‌توان با ساخت مؤلفه به شیوه‌ای که بعداً خواهد آمد، از این خصوصیات برای بهبود دقت عملیات ترازبندی استفاده نمود.

به طور کلی نوآوری‌های موجود در این روش را می‌توان در دو بخش طبقه بندی کرد.

در بخش اول و در قسمت استخراج ویژگی، استفاده از مولفه های زبانی برچسب دستوری در ساخت مؤلفه های بردار ویژگی به عنوان یک رویکرد جدید مطرح می‌باشد. همچنین در همین قسمت استفاده از اصلاحگر های برچسب زن دستوری نیز از دیگر مواردی است که تا کنون به آن پرداخته نشده است. همچنین بکارگیری ابزار مترجم ماشینی گوگل به عنوان یک ابزار کمکی سبب شده است که عملیات ساخت ترجمه تحت‌اللفظی عبارات انگلیسی، ارزش دقیق‌تری را به میزان نزدیکی دو عبارت فارسی و انگلیسی کاندید برای ترازبندی تولید نماید.

³ Decision tree

⁴ Support Vector Machine