


	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیرپروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	



عنوان زیرپروژه:

مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی

	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
	مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی	کد زیر پروژه: پیکرمتن فارس - 3 - ث	ویرایش: 1/0	تاریخ: 1388/04/19

فهرست مطالب

شماره صفحه	عنوان
3.....	1. مقدمه
5.....	2. مروری بر کیفیت پیکره های مخلوط فارسی و عربی
8.....	3. رده بندی و تشخیص زبان متن
11.....	4. تشخیص زبان در متون تک-زبانه
12.....	5. تشخیص زبان در متون چند زبانه و الگوریتم Ludovik and Zacharski
13.....	1-5. مرحله پیش پردازش و آموزش الگوریتم Ludovik and Zacharski
16.....	2-5. مرحله رده بندی و تشخیص الگوریتم Ludovik and Zacharski
17.....	3-5. مرحله تقطیع الگوریتم Ludovik and Zacharski
20.....	4-5. ارزیابی الگوریتم Ludovik and Zacharski
21.....	6. نتیجه گیری
22.....	7. پایگاه های on-line برای تشخیص زبان متن
23.....	مراجع

	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		کد زیر پروژه: پیکرمتن فارس - 3 - ث	ویرایش: 1/0	تاریخ: 1388/04/19

1. مقدمه



به جهت امتزاج فرهنگ فارسی با فرهنگ اسلامی، متون نظم و نثر فارسی مملو از جملات و عبارتهای عربی است. بخش عمده این قبیل متون متعلق به حوزه علوم اسلامی شامل علوم قرآن و حدیث و نیز فقه است.

قدر مسلم پردازش رایانه‌ای چنین متونی مستلزم سامانه‌های جداگانه برای فارسی و عربی است که ایجاد می‌کند سامانه‌ای مستقل ابتدا بخش‌های فارسی را از عربی تفکیک نماید. برای دستیابی به راه‌کاری مناسب برای این مهم طرح پژوهشی حاضر حصول به پاسخ سؤالات ذیل را سرلوحه پژوهش امکان‌سنجی خود قرار داده است تا روش‌های تشخیص هوشمند جمله فارسی از عربی در پیکره‌های مخلوط فارسی و عربی تبیین شده و در پژوهش‌های آتی به سامانه‌های هوشمند تبدیل شوند.

1. روش‌های کلی تشخیص هوشمند جمله‌های یک زبان در پیکره‌های مخلوط این زبان با زبان‌های مشابه چیست؟
2. کیفیت و کمیت پیکره‌های متنی مخلوط فارسی و عربی برای تشخیص هوشمند جمله فارسی از عربی در این پیکره‌ها چگونه باید باشد؟
3. چه نرم افزارهایی برای پردازش پیکره‌های متنی مخلوط فارسی و عربی برای تشخیص هوشمند جمله فارسی از عربی در این پیکره‌ها مورد نیاز است؟
4. در صورت نیاز به پایگاه دانش برای تشخیص هوشمند جمله فارسی از عربی در پیکره‌های متنی مخلوط فارسی و عربی، آیا این پایگاه دانش مبتنی بر قواعد پیاده‌سازی شده انسانی خواهد بود یا آموخته‌های ماشین از پیکره‌های متنی مخلوط فارسی و عربی؟
5. در پیکره‌های متنی مخلوط فارسی و عربی چه ویژگی‌هایی وجه ممیزه بین جمله‌های فارسی و عربی است و بهترین شیوه‌گزینش این ویژگی‌ها چیست؟

سؤالات فوق پژوهش حاضر را از بررسی‌های ذیل به پاسخ‌های لازم سوق داده است.



1. روش‌های کلی رده بندی متون،

	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیرپروژه:			
مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		کد زیرپروژه: پیکرمتن فارس - 3 - ث	ویرایش: 1/0	تاریخ: 1388/04/19

2. سامانه‌هایی در دنیا که توانسته‌اند در داخل متنی که تمام آن با کاراکترهای یکسان نگارش شده، رشته‌ای را که به زبانی به غیر از زبان کلی متن است تشخیص دهند،
3. کیفیت و کمیت پیکره‌های متنی مخلوط فارسی و عربی،
4. کیفیت و کمیت بخش‌های عربی، از قبیل آیات قرآن کریم، احادیث و اشعار عربی، در داخل پیکره‌های متنی مخلوط فارسی و عربی،
5. نرم افزارهای پردازش دستوری پیکره‌های متنی فارسی و عربی،
6. پژوهش‌ها و سامانه‌هایی که با هدف پردازش و طبقه بندی متون عربی وجود دارند،
7. پژوهش‌ها و سامانه‌هایی که با هدف پردازش و طبقه بندی متون فارسی وجود دارند،
8. پژوهش‌ها و سامانه‌هایی که با هدف پردازش و طبقه بندی پیکره‌های متنی مخلوط فارسی، عربی و زبان‌های مشابه دیگر وجود دارند.

در گزارش پژوهشی حاضر ابتدا کیفیت پیکره‌های مخلوط فارسی و عربی مورد بررسی قرار می‌گیرد. آنگاه مشکل تشخیص هوشمند جمله فارسی از عربی در پیکره‌های مخلوط فارسی و عربی از لحاظ علمی تبیین شده و شاخه علمی که برای حل این مسئله از آن مدد جسته خواهد شد، یعنی رده بندی متون، توضیح داده می‌شود. به عنوان نمونه کاربرد این نگرش علمی، به سامانه‌ای که بطور اولیه به این منظور با همکاری نگارنده این گزارش پژوهشی و گروهی از دانشجویان در دانشگاه نبی اکرم (ص)، تبریز، پیاده سازی شده است اشاره خواهد گردید. سپس در همین راستا، به عنوان نمونه بهینه که در طی دهه گذشته جامع ترین الگوریتم عرضه عمومی شده است، الگوریتم¹ Ludovik and Zacharski بطور کامل تشریح می‌شود. در پایان برخی وب‌گاه‌های عمومی که در همین راستا خدمات on-line ارائه می‌کنند معرفی خواهند شد.

¹ Ludovik and Zacharski

	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیرپروژه:			
مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		کد زیرپروژه: پیکرمتن فارسی - 3 - ث	ویرایش: 1/0	تاریخ: 1388/04/19

2. مروری بر کیفیت پیکره‌های مخلوط فارسی و عربی

در بررسی کیفیت و کمیت امتزاج متون فارسی و عربی در یک دیگر آنچه در پژوهش حاضر بیشتر مورد توجه قرار گرفته است متونی است که عمده آنها به زبان فارسی‌اند و تنها بخش‌هایی از آنها به زبان عربی است. اینگونه متون بطور غالب متون اسلامی یا ادبی‌اند که در آنها گونه عربی استفاده شده در بخش‌های عربی شامل آیاتی از قرآن کریم، احادیث و روایات، و نیز اشعار عربی است. در این میان، تفاسیر فارسی و یا ترجمه شده از عربی برای قرآن کریم بهترین نمونه‌ای هستند که، در عین حجم بسیار بالا، از تمامی گونه‌های عربی فوق در آنها یافت می‌شود. بنابراین در پژوهش حاضر این متون به عنوان نمونه اصلی مورد بررسی قرار گرفته‌اند.

نمونه متنی که در پی می‌آید بخشی کوچک از کتاب تفسیر نور در باره سوره یونس است.

دهمین سوره‌ی قرآن کریم که در اوایل بعثت در مکه نازل شده است، «یونس» نام دارد. این سوره یکصد و نه آیه دارد و عمده‌ی مطالب آن، پیرامون توحید و حقایق قرآن، پاسخ به منکران وحی، بیم دادن مشرکان، بیان عظمت آفرینش و آفریدگار، ناپایداری دنیا و توجه دادن به آخرت است.



...

أَكَانَ لِلنَّاسِ عَجَبًا أَنْ أَوْحَيْنَا إِلَى رَجُلٍ مِنْهُمْ أَنْ أَنْذِرِ النَّاسَ وَبَشِّرِ الَّذِينَ آمَنُوا أَنْ لَهُمْ قَدَمَ صِدْقٍ عِنْدَ رَبِّهِمْ قَالَ الْكَافِرُونَ إِنَّ هَذَا لَسَاحِرٌ مُبِينٌ (2)

آیا برای مردم شگفت‌آور است که به مردی از خود آنان وحی کردیم که مردم را بیم و هشدار بدهد و به مؤمنان بشارت بدهد که برای آنان نزد پروردگارشان جایگاه نیکویی است؟ کافران گفتند: همانا این مرد جادوگری آشکار است!

سوره‌ی قبل، یعنی توبه، رفتار منافقان و کيفر آنها را بیان می‌کند این سوره به بیان رفتار مشرکان پرداخته است.

...

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیرپروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	

پیام ها:

1 ریشه‌ی کفر، غالباً استبعاد و تعجب از وحی است. «عَجَبًا أَنْ أَوْحَيْنَا» 2 سرچشمه‌ی دعوت انبیا، وحی الهی است. «أَوْحَيْنَا إِلَى رَجُلٍ» 3 لیاقت‌های معنوی افراد، در ظاهر دیده نمی‌شود. اگر کسی مورد لطف خاص خدا قرار گرفت، او را تحمل کنیم. «رَجُلٍ مِنْهُمْ» 4 پیامبران به خاطر الگو بودنشان باید از مردم و درد آشنا باشند. «مِنْهُمْ» 5 وظیفه‌ی پیامبران، بشارت و انذار است. «أَنْذِرِ»، «بَشِّرِ» 6 هر آنچه نمی‌فهمیم، آن را رد نکنیم. زیرا که انکار و تهمت زدن، شیوه‌ی کفار است. «قَالَ الْكَافِرُونَ إِنَّ هَذَا لَسَاحِرٌ مُّبِينٌ» 7 ایمان، زمینه‌ی قدم صدق و جایگاه ویژه نزد خداوند است. «بِأَنَّ لَهُمُ الْجَنَّةَ»

همان طور که مشاهده می‌شود در نیمه اول این متن آیه مورد نظر به طور مجزا در بین دو بند توضیح فارسی آمده است؛ ولی در ادامه متن همان آیه شکسته شده و به صورت عبارتها و یا کلماتی در دل توضیحات فارسی ظاهر شده است. این به آن معناست که، در پردازش این متن و جدا کردن بخش‌های عربی از فارسی، واحد پردازش می‌تواند جمله کامل، عبارت و یا حتی کلمه تنها باشد.



البته باید توجه داشت در این متن بخش‌های عربی استفاده شده نه تنها دارای حرکه می‌باشند، بلکه به هنگام آماده سازی متن به طور دستی از بخش‌های فارسی جدا شده اند. این شرایط همیشه برقرار نیست. به عنوان مثال در کتاب تفسیر نمونه از قرآن کریم در خصوص همان سوره که بخشی از آن در پی می‌آید در ابتدا و انتهای متن از عبارتها و جمله‌های عربی استفاده شده است که مشخصه‌های فوق را ندارند.

سوره یونس (ع) این سوره در مکه نازل شده و 109 آیه است

بسم الله الرحمن الرحيم

محتوی و فضیلت این سوره

این سوره که از سوره‌های مکی است، و به گفته بعضی از مفسران بعد از سوره اسراء و قبل از سوره هود نازل شده است همانند بسیاری از سوره‌های مکی روی

	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		کد زیر پروژه: پیک متن فارس - 3 - ث	ویرایش: 1/0	تاریخ: 1388/04/19

چند مساله اصولی و زیر بنائی تکیه می کند، که از همه مهمتر مساله " مبدء " و " معاد " است.

منتها نخست از مساله وحی و مقام پیامبر ص سخن می گوید، سپس به نشانه هایی از عظمت آفرینش که نشانه عظمت خدا است می پردازد، بعد، مردم را به ناپایداری زندگی مادی دنیا و لزوم توجه به سرای آخرت و آمادگی برای آن از طریق ایمان و عمل صالح متوجه می سازد.

و به تناسب همین مسائل قسمت های مختلفی از زندگی پیامبران بزرگ از جمله نوح و موسی و یونس (ع) را بازگو می کند و به همین مناسبت نام سوره یونس بر آن گذارده شده است.



و باز برای تایید مباحث فوق، سخن از لجاجت و سرسختی بت پرستان به میان می آورد، و حضور و شهود خدا را در همه جا برای آنها ترسیم می کند، و مخصوصا برای اثبات این مساله از اعماق فطرت آنان که به هنگام مشکلات آشکار می شود و به یاد خدای واحد یکتا می افتند، کمک می گیرد.

و بالاخره برای تکمیل بحثهای فوق در هر مورد مناسبی از بشارت و انذار، بشارت به نعمتهای بی پایان الهی برای صالحان و انذار و بیم دادن طاغیان و گردنکشان، استفاده می کند.

لذا در روایتی از امام صادق ع می خوانیم که فرمود:

من قرء سورة یونس فی کل شهرین او ثلاثه لم یخف علیه ان یكون من الجاهلین و کان یوم القیامة من المقربین

نکته قابل توجه دیگر اینکه در نمونه متن های ارائه شده، از لحاظ کمی، بخش های عربی سهم کمی را به عهده دارند.



	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	

3. رده بندی و تشخیص زبان متن

مسئله تشخیص و تفکیک بخش‌های تک-زبانۀ در داخل متون چند-زبانۀ بطور رایج به عنوان یک مسئله تشخیص زبان^۱ شناخته شده و از منظر رده بندی متن^۲ نگریسته می‌شود. بطور کلی رده بندی متن، که انتساب اسناد متنی بر اساس محتوی به یک یا چند طبقه از قبل تعیین شده است، یکی از مهم‌ترین مسایل در متن کاوی است. مرتب کردن بلادرنگ نامه‌های الکترونیکی یا فایل‌ها در سلسله مراتبی از پوشه‌ها، تشخیص موضوع متن، جستجوی ساخت یافته و یا پیدا کردن اسنادی که در راستای علایق کاربر می‌باشد، از جمله کاربردهای مبحث رده بندی (طبقه بندی، دسته بندی یا کلاس بندی) متن است. از این نظر بطور کلی در متون چند زبانۀ^۳، سامانه تشخیص زبان^۴ باید هر قطعه تک زبانۀ^۵ و نیز زبان هر قطعه را مشخص نماید. در ایجاد سامانه‌های ترجمه ماشینی^۶ اغلب پیکره‌های^۷ بسیار بزرگی مورد نیاز است و برای گردآوری خودکار چنین پیکره‌ای می‌توان یک سامانه تشخیص زبان را با یک عنکبوت وب^۸ ترکیب کرد. الگوریتم‌های تشخیص نه تنها زبان متن الکترونیکی بلکه encoding کاراکترهای زبان آن را تشخیص می‌دهند.

پس می‌توان مسئله تشخیص زبان قطعه‌ای از یک متن را به عنوان مسئله رده بندی متن دانست که در آن سه مرحله پیش-پردازش^۹، آموزش رده بند^{۱۰} و رده بندی^{۱۱} وجود دارد.

-
- Language detection^۱
 - Classification text^۲
 - Multilingual documents^۳
 - Language recognition system^۴
 - Monolingual segment^۵
 - Machine translation^۶
 - Corpora^۷
 - Web spider^۸
 - Pre-processing^۹
 - Training classifier^{۱۰}
 - Classification^{۱۱}

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	

در مرحله پیش-پردازش، دانش موجود در هر متن باید بازنمایی^۱ شود تا قابل استفاده نرم افزارهای رده بندی گردد. این بازنمایی به شکل مدل برداری^۲ از ویژگی‌های^۳ متن که برگرفته از عناصر موجود در آن است انجام می‌گیرد. رایج‌ترین این ویژگی‌ها کلمه^۴ است که برخی مواقع به وسیله نرم افزارهای پردازش زبان طبیعی^۵ با برجسب‌هایی^۶ همراه می‌شوند که حاوی اطلاعات صرفی-نحوی^۷ آنها است. در برخی موارد این ویژگی‌ها به صورت‌های پیچیده‌تری نیز تبدیل می‌شوند که رایج‌ترین شیوه استفاده از مدل-سازی چند-گرمی^۸ است. برای بازنمایی متون با استفاده از ویژگی‌های مناسب استخراج یا ساخته شده از همان متون هر یک از این ویژگی‌ها به ارزشی عددی نگاشت داده می‌شود که به عنوان وزن^۹ آن ویژگی در متن مورد نظر محاسبه می‌شود. برای این منظور ماتریس $m \times n$ ایجاد می‌شود که در آن m کل متن-های موجود در رده‌ها، n کل ویژگی‌های ایجاد شده، و A_{ij} تعداد تکرار ویژگی i یا به عبارتی وزن آن در متن j است.

در این هنگام بطور معمول مشکل گستره و پراکندگی زیاد مقدار-وزن این ویژگی‌ها پیش می‌آید که برای حل آن به شیوه‌های گزینش ویژگی‌ها^{۱۰} رجوع می‌شود تا مقادیری از ویژگی‌ها که بیش از دیگران قابلیت تمییز دهندگی^{۱۱} متون را دارند گلچین شوند.

در مرحله آموزش رده بند، برای ایجاد سامانه‌های رده بندی از سامانه‌های یادگیری ماشین با ناظر^{۱۲} بهره‌گیری می‌شود. این سامانه‌ها بوسیله مقدار-وزن‌های به دست آمده در مرحله پیش-پردازش از متونی

^۱ Knowledge representation

^۲ Model vector

^۳ Features

^۴ Word

^۵ Natural language processing

^۶ Tags

^۷ Morpho syntactic



^۸ N-gram

^۹ Weight

^{۱۰} Selection feature



^{۱۱} Discrimination

^{۱۲} Supervised machine learning

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیرپروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	

که تحت عنوان متون آموزشی^۱ از قبل رده بندی شده‌اند آموزش داده شده و به یک رده بند متون^۲ تبدیل می‌شوند.

در مرحله رده بندی، متون مورد نظر برای رده بندی، پس از گذر از مرحله پیش-پردازش و تبدیل به بردارهای مقدار-وزن در قالب ویژگی‌های انتخاب شده در مرحله آموزش، به رده بند داده می‌شوند تا، مطابق با رده‌های از قبل آموزش داده شده به رده بند، در یکی از رده‌های آن رده بندی شوند. میزان موفقیت رده بند در این رده بندی معادل ارزیابی انجام شده بر روی آن در مرحله آموزش فرض می‌شود.



	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		کد زیر پروژه: پیک متن فارس - 3 - ث	ویرایش: 1/0	تاریخ: 1388/04/19

4. تشخیص زبان در متون تک-زبانه

مطابق گزارش اخیر [3]، اصلی ترین تفاوت بین شیوه های متداول تشخیص زبان در متون تک-زبانه وابسته به مدل زبانی است که از متن مورد نظر می سازند (هرچند در مرحله آموزش رده بند در بکارگیری الگوریتم های مختلف یادگیری ماشین با هم تفاوت هایی دارند). بر این اساس این شیوه ها از سه نگرش کلی پیروی می کنند که وابسته به سه نوع ویژگی بکار رفته در مرحله پیش-پردازش اند:

1. کلمات کوتاه، که در آن مدل زبانی فقط از کلمات کوتاه با طول محدود و بدون توجه به تواتر این کلمات ساخته می شود. به طور نمونه، [4] از تمامی کلمات با طول حداکثر پنج کاراکتر که حداقل سه بار تکرار شده باشند بهره می برد و بدین ترتیب مدل زبانی اش از 980 تا 2750 چنین کلمه هایی استفاده می کند.
2. کلمات متواتر، که در آن مدل زبانی با استفاده از کلماتی ساخته می شود که بین کلیه کلمات متن بیشترین تواتر را دارند. به طور نمونه، [5] با استفاده از یک صد کلمه بایشترین تواتر جداولی می سازد که در آنها هر کلمه ارزش تواتری خاصی دارد که خارج قسمت تقسیم تواتر نسبی هر کلمه در جدول بر تواتر نسبی کلمه اول آن جدول است.
3. چند-گرم کلمات، که در آن توالی کلیه کلمات متن در نظر گرفته شده و از آنها بر اساس چند-گرمی های با طول ثابت یا متغیر مدل زبانی ساخته می شود.

شیوه غالب کنونی در انتخاب ویژگی های تمییز دهنده برای شناسایی زبان استفاده از شیوه سوم است که به عنوان مثال در الگوریتم پیشنهادی [6] بکار رفته است و با ایجاد مدل زبانی با 300 چند-گرمی نتیجه 99.8% به دست آورده است. با این وجود و "در عین سادگی و شهودی بودن و نیز موفقیت در بسیاری موارد، این شیوه چندین نقطه ضعف دارد. اول اینکه، علیرغم کارکرد خوب در متون با اندازه متوسط، با کاهش طول متن عملکرد آن کاهش چشم گیری می یابد. دوم اینکه، به جهت اتکا بر واژه ها، این شیوه وابستگی بسیاری به توانایی سامانه در تقطیع متن به واژه ها دارد" [7]. ولی مهمترین کاستی تمامی این شیوه ها اختصاص آنها به متون تک-زبانه است و این در حالی است که پیکره مورد هدف این پژوهش بطور کلی پیکره ای دو زبانه حاوی فارسی و عربی است. بنابراین الگوریتم هایی در تشخیص زبان متون باید بیشتر مد نظر قرار گیرند که برای متون چند-زبانه طراحی شده اند.

	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		کد زیر پروژه: پیک-متن فارس - 3 - ث	ویرایش: 1/0	تاریخ: 1388/04/19

5. تشخیص زبان در متون چند زبانه و الگوریتم Ludovik and Zacharski

در بررسی الگوریتم‌های بکار رفته برای تشخیص زبان در متون چند-زبانه نکته اولی که جلب توجه می‌کند نادر بودن چنین الگوریتم‌هایی است. نکته قابل توجه اصلی در این الگوریتم‌ها، از نمونه‌های قبلی مانند [8] گرفته تا الگوریتم جدید استفاده شده در [9] همگی برای متون حاوی مخلوطی از زبان‌هایی به غیر از فارسی و عربی طراحی شده‌اند. در تلاشی برای حل کاستی موجود و با هدف تشخیص زبان در پیکره‌های مخلوط چند-زبانه که فارسی و عربی را نیز در بر می‌گیرد، به شیوه‌ای مشابه آنچه که در شیوه اول انتخاب ویژگی‌های متنی در تشخیص زبان متون تک-زبانه مورد استفاده است سامانه‌ای در دانشگاه نبی اکرم (ص) در تبریز به صورت نمونه‌ی اولیه^۱ طراحی و پیاده سازی شده است [10]. مشابه موفقیت‌های گزارش شده برای این شیوه در تشخیص زبان متون تک-زبانه، این سامانه در تمییز جملات فارسی و عربی از یکدیگر توفیق زیادی داشته است. علیرغم این موفقیت نسبی دارای کاستی مهمی است. این سامانه فقط موقعی موفقیت دارد که متن بطور کامل در واحد جمله تقطیع شده باشد؛ بنابراین سامانه فوق در تشخیص وا حد‌های کوچکتر از جمله همانند عبارت یا کلمه هنوز توفیقی ندارد. و این در حالی است که، مطابق بخش 2 این گزارش پژوهشی، در پیکره‌های مخلوط فارسی و عربی نه تنها از جمله‌های عربی بلکه از واحدهای کوچکتری چون عبارت و کلمه‌های عربی نیز استفاده شده است.

بر پایه این توضیحات، الگوریتمی مورد نیاز است که بتواند نه تنها در تشخیص طول رشته‌های بلند مانند جمله و کوتاه مانند عبارت موفق عمل کند بلکه در متون مخلوط حاوی فارسی و عربی نیز امتحان خود را پس داده باشد. الگوریتم Ludovik and Zacharski در [7] چنین الگوریتمی است و در واقع از این جهت تنها الگوریتم گزارش شده و قابل دسترس است. Ludovik and Zacharski در الگوریتم خود برای شناسایی زبان متون چند زبانه از ترکیب چند-گرم‌ها^۲ با مرتبه‌های مختلط^۳، زنجیره مارکف^۴، درست‌نمایی حداکثر^۵ و برنامه‌سازی پویا^۱ استفاده کرده‌اند که وابستگی این شیوه را به تشخیص کلمات



^۱ Prototype

^۲ N-gram

^۳ Mixed order

^۴ Markov chain

^۵ Maximum likelihood

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	

متن می‌کاهد. این الگوریتم دارای دو مرحله اصلی آموزش و رده بندی و نیز یک مرحله فرعی تقطیع^۲ است.

5-1. مرحله پیش پردازش و آموزش الگوریتم Ludovik and Zacharski

در مرحله پیش پردازش و آموزش، الگوریتم Ludovik and Zacharski، برای ایجاد ویژگی‌ها و مقدار-وزن آنها، مجموعه‌ای خالی را مورد استفاده قرار می‌دهد که از چند-گرمی‌های با طول متغیر برگرفته شده از هر متن آموزشی پر خواهد شد (در پژوهش گزارش شده تا چهار-گرمی استفاده شده است). ابتدا فهرست‌های جداگانه‌ای برای تک-گرمی‌ها، جفت-گرمی‌ها و غیره ایجاد می‌شود. آنگاه بخشی از این چند-گرمی‌ها گزینش شده و به مجموعه اولیه اضافه می‌گردد - ابتدا تک-گرمی‌ها، سپس جفت-گرمی‌ها و به همین ترتیب. گزینش با محاسبه وزن آموزشی^۳ انجام می‌گیرد که بطور مثال برای تک-گرمی‌ها با استفاده از (1) است.

$$W(a_1) = -p(a_1) \log p(a_1) \quad (1)$$



سپس برای محاسبه وزن آموزشی هر k-گرمی $(a_1 a_2 \dots a_k)$ فرایند بازگشتی ذیل به ازای $k = 2 \dots N$ طی می‌شود. اگر (k-1)-گرمی $(a_2 \dots a_k)$ هنوز در مجموعه اولیه نیست وزن آموزشی با استفاده از (2) محاسبه می‌شود و در غیر این صورت از (3).

$$W(a_1 a_2 \dots a_k) = -p(a_1 a_2 \dots a_k) \log p(a_k | a_1 a_2 \dots a_{k-1}) \quad (2)$$

^۱ Dynamic programming

^۲ Segmentation

^۳ Training weight

	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
	مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی	کد زیر پروژه: پیکرمتن فارس - 3 - ث	ویرایش: 1/0	تاریخ: 1388/04/19

$$W(a_1 a_2 \dots a_k) = -p(a_1 a_2 \dots a_k) (\log p(a_k | a_1 a_2 \dots a_{k-1}) - \log p(a_k | a_2 a_3 \dots a_{k-1})) \quad (3)$$

سپس با مرتب‌سازی نزولی تمامی چند-گرمی‌ها بر اساس وزن‌های آموزشی، بالاترین آنها انتخاب و در مجموعه اولیه قرار می‌گیرند. در انتهای این فرایند، مجموعه اولیه شامل اجتماع چند-گرمی‌های گزینش شده برای همه L زبان‌ها در مجموعه آموزشی است.

آنگاه به ازای هر چند-گرمی ($n = 1 \dots N$) در مجموعه اولیه بردار L-بعدي وزن تشخیص اولیه¹ به ترتیب ذیل محاسبه می‌شود. هر i-امین عضو بردار، PRW، مطابق i-امین زبان بوده و شامل یک وزن تشخیص است. وزن تشخیص برای تک-گرمی‌ها به صورت (4) و برای تمامی دیگر چند-گرمی‌ها به صورت (5) خواهد بود.



$$PRW_i(a_1) = -\log p_i(a_1) \quad (4)$$

$$PRW_i(a_1 a_2 \dots a_k) = -\log p_i(a_k | a_1 a_2 \dots a_{k-1}) \quad (5)$$

در واقع (5) یاد-لگاریتم احتمال شرطی a_k است به شرط بافتار $a_1 a_2 \dots a_{k-1}$ در مجموعه آموزشی مطابق با i-امین زبان. اگر یک چند-گرمی در مجموعه اولیه در داده آموزشی زبانی پیش نیاید، وزن تشخیص آن برای این زبان بعنوان ارزش حداکثری MAX تعریف می‌شود. در پایان محاسبه، بردار L-بعدي وزن تشخیص هر چند-گرمی بصورت (6) محاسبه می‌شود.

$$RW(a_1 a_2 \dots a_N) = \begin{cases} PRW(a_1 a_2 \dots a_N) & | a_1 a_2 \dots a_N \in CNP \\ PRW(a_2 a_3 \dots a_N) & | a_2 a_3 \dots a_N \in CNP \\ \dots & | \dots \\ PRW(a_N) & | a_N \in CNP \end{cases} \quad (6)$$

¹ PRW: primary recognition weight

	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
	مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی	کد زیر پروژه: پیکرمتن فارس - 3 - ث	ویرایش: 1/0	تاریخ: 1388/04/19

vector of MAX | otherwise

در پایان، میانگین وزن^۱ و پراکندگی وزن^۲ برای هر زبان محاسبه می شود که به ترتیب ذیل انجام می - گیرد. هر متن آموزشی زبان فرضی i به K قطعه 500 بایتی $(x_1x_2...x_{500})$ تقطیع می شود. آنگاه به ازای هر قطعه k (7) محاسبه می شود.

$$d_{i,k} = \sum_{s=1}^{500} \frac{RW_i(x_{s-N+i}x_{s-N+2}...x_s)}{500} \quad (7)$$



میانگین وزن برای زبان i بر اساس (8) تعریف می شود.

$$WA_i = \frac{\sum_{k=1}^K d_{i,k}}{K} \quad (8)$$

و پراکندگی برای زبان i بر اساس (9) تعریف می شود.

$$D_i^2 = \frac{\sum_{k=1}^K (d_{i,k} - WA_i)^2}{K} \quad (9)$$

^۱ Weight average
^۲ Weight dispersion

	عنوان پروژه:			
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی			
	عنوان زیر پروژه:			
	مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی	کد زیر پروژه: پیکرمتن فارس - 3 - ث	ویرایش: 1/0	تاریخ: 1388/04/19

2-5. مرحله رده بندی و تشخیص الگوریتم Ludovik and Zacharski

مرحله تشخیص شامل دو فرایند است: رده بندی و تأیید. در رده بندی متن تحت تشخیص به طور اولیه در رده یکی از زبان های تعیین شده در مرحله آموزش رده بندی می شود. آنگاه در فرایند تأیید اینکه متن مورد نظر با زبان پیشنهاد شده در فرایند اول چقدر مناسب دارد تعیین می شود؛ اگر تناسب کافی نباشد، زبان متن به عنوان ناشناس رده بندی می شود.

فرایند رده بندی: فرض می شود $x_0 x_1 \dots x_{s-1}$ رشته بایستی تحت رده بندی باشد. یک بردار وزن نتیجه تشخیص¹ به صورت (10) تعریف می شود. نتیجه تشخیص (زبان پیشنهادی برای متن) به صورت (11) خواهد بود (زبان متن به عنوان زبان i^* رده بندی می شود که متناظر است با آن عضوی از بردار وزن نتیجه تشخیص که کمترین مقدار را دارد).



$$RRW = \frac{\sum_{s=0}^{S-1} RW(x_{s-N+1} x_{s-N+2} \dots x_s)}{S} \quad (10)$$

$$i^* = \arg_s \min RRW_i \quad (11)$$

فرایند تأیید: اگر رابطه (12) برقرار باشد، آنگاه زبان متن به عنوان اینکه زبان آن i^* است رده بندی می شود. اما اگر این رابطه برقرار نباشد، آنگاه زبان آن به عنوان ناشناس رده بندی خواهد شد.

$$\frac{RRW_{i^*} - WA_{i^*}}{D_{i^*}} \leq VER_THR \quad (12)$$

¹ Result Recognition Weight

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	

اصلی ترین شرط موفقیت الگوریتم Ludovik and Zacharski در مرحله رده بندی و تشخیص این است که متن مورد نظر تک زبانه باشد. برآورده شدن این شرط به عهده مرحله تقطیع است.

3-5. مرحله تقطیع الگوریتم Ludovik and Zacharski



از آنجا که متون تحت بررسی چند زبانه هستند، این متون باید که ابتدا به قطعات تک زبانه تقطیع شده و سپس زبان هر قطعه تشخیص داده شود. این امر با افزودن یک الگوریتم برنامه نویسی پویا مبتنی بر مدل مارکوفی ساده از متون چند زبانه امکان پذیر است. این الگوریتم، که اولین بار توسط Vintsiuk برای تقطیع استفاده شد، گونه ای بسط یافته از الگوریتمی است که کارآیی محاسباتی آن توسط Ludovik افزایش یافته است و دارای دو مرحله است.

1. مرحله تشخیص، که الگوریتم برنامه نویسی پویا را بکار می گیرد تا با به حداکثر رسانی احتمال رشته کلی کاراکترهای متن تحت پردازش به تقطیع دست یابد،
2. مرحله تأیید، که تعیین می کند تا چه حدی هر قطعه با زبان تخصیص یافته مناسب دارد.

در مرحله تشخیص، الگوریتم تقطیع بر مدل مارکوفی ذیل از متن چند زبانه بنا شده است. این مدل به ازای هر زبان یک حالت 1 دارد، $1 \leq i \leq L$ ، بعلاوه یک حالت 0 اضافی برای قطعاتی که در هیچ یک از زبان های آموزشی سامانه وجود ندارند، از جمله قطعاتی که به هیچ عنوان به زبانی متعلق نیست (مانند جدولی از اعداد). اگر سامانه در حالت i باشد، بسته به بافتار 2 و بر اساس مدل چند-گرمی، کاراکترهای زبان i را تولید خواهد کرد. تغییر وضعیت از زبان i_1 به زبان i_2 بر اساس احتمالات انتقال $p(i_2|i_1)$ و توزیع احتمال طول قطعه $p(r)$ ، $r_{\min} \leq r \leq r_{\max}$ ، با گستره تغییر بین طول قطعه حداقل r_{\min} تا حداکثر r_{\max} تعریف می شود. الگوریتمی که در ذیل ارائه می شود با توجه به مدل مارکوفی نسبت به حداکثر درستنمایی متن مورد مشاهده حصول اطمینان کرده، قطعه ها و زبان های قطعه ها را می یابد.

¹ State

² Context

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیرپروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	

اگر k امین قطعه در s_{k-1} شروع شده و در s_k-1 تمام شود، آنگاه با استفاده از پادلگاریتم، معیار حداقل سازی بر اساس (13) است که در آن TSW یک وزن قطعه است مطابق (14) به شرط صفر نبودن i_k ، وگرنه مطابق (15) (ضریب ثابت JUNK_THR)، به این منظور به کار می رود تا قطعه های متنی را که در هیچ یک از زبان های آموزشی نیستند حذف نماید):



$$Q(\{x_s, 0 \leq s \leq S\}, \{s_k, i_k, 0 \leq k \leq K\}) = \sum_{k=1}^K TSW(i_{k-1}, i_k, s_{k-1}, s_k) \quad (13)$$

$$TSW(i_{k-1}, i_k, s_{k-1}, s_k) = \sum_{s=s_{k-1}}^{s_k-1} RW_{i_k}(x_{s-N+1} x_{s-N+2} \dots x_s) - \log p(i_k | i_{k-1}) - \log p(s_k - s_{k-1}) \quad (14)$$

$$TSW(i_{k-1}, i_k, s_{k-1}, s_k) = \text{JUNK_THR}^*(s_k - s_{k-1}) - \min_{i:1 \leq i \leq L} [\sum_{s=s_{k-1}}^{s_k-1} RW_i(x_{s-N+1} \dots x_s) - \log p(i_k | i_{k-1})] - \log p(s_k - s_{k-1}) \quad (15)$$

الگوریتم برنامه نویسی پویا که مقادیر بهینه $\{0 < k < K, i_k, s_k\}$ را می یابد شامل گام تکراری (17) برای (16) با مقادیر اولیه (18) می شود.

Junk threshold¹

	عنوان پروژه:		
	فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیرپروژه:		
	مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی	کد زیرپروژه: پیک-متن فارسی - 3 - ث	ویرایش: 1/0
	تاریخ: 1388/04/19		

$$\begin{aligned}
 Ind(i_k, s_k) = (i_{k-1} * , s_{k-1} *) = & \arg \min_{(i_{k-1}, s_{k-1});} (F(s_{k-1}, i_{k-1})) \\
 & i_{k-1}; 0 \leq i_{k-1} \leq L; \\
 & s_{k-1}; r_{\min} \leq (s_k - s_{k-1}) \leq r_{\max}; s_{k-1} \geq 0 \\
 + TSW(i_{k-1}, i_k, s_{k-1}, s_k) & \quad (16)
 \end{aligned}$$

$$\begin{aligned}
 F(i_k, s_k) = F(i_{k-1} * , s_{k-1} *) \\
 + TSW(i_{k-1} * , i_k, s_{k-1} * , s_k) \quad (17)
 \end{aligned}$$



$$F(i, 0) = 0, 0 \leq i \leq L, F(i, s) = \text{INFINTY}, s: 0 < s < r_{\min}, 0 \leq i \leq L \quad (18)$$

بعد از اینکه الگوریتم گامهای بازگشتی را تمام می کند، زبان آخرین قطعه در مجموعه بهینه قطعات (19) یافت می شود که در آن S طول کل مربوط به متن است.

$$i * = \arg \min_{i: 0 \leq i \leq 1=L} F(i, S) \quad (19)$$

مقدار i^* به همراه اطلاعات ذخیره شده در آرایه $Ind(i, s)$ اجازه خواهد داد تا تمامی قطعه ها با برچسبهای زبان مربوطه از مجموعه بهینه دریافت شوند. در طی این فرایند، تعداد بهینه قطعه ها به طور خودکار تعیین می شود.



فرایند تأیید در مرحله تقطیع همانند فرایند تأیید در مرحله تشخیص و رده بندی متون تک زبانه است.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	

4-5. ارزیابی الگوریتم Ludovik and Zacharski



برای اندازه‌گیری نرخ خطای این الگوریتم در تقطیع متون چند-زبانه، الگوریتم روی 6 متن که مخلوطی از 1000 قطعه با طول‌های به طور نسبی مساوی (بین 20 تا 1000 بایت) از زبان‌های مختلف از جمله عربی و فارسی بودند توسط طراحانش ارزیابی شده است. نرخ خطای این الگوریتم 0.47% برای قطعه متن‌های 1000 بایتی، 0.69% برای قطعه متن‌های 540 بایتی، 1.4% برای قطعه متن‌های 202 بایتی، 2.08% برای قطعه متن‌های 101 بایتی، 4.7% برای قطعه متن‌های 49 بایتی، و 12.88% برای قطعه متن‌های 20 بایتی (حدود سه کلمه) گزارش شده است و این نرخ آخری بدین معناست که الگوریتم فقط در تشخیص بایت‌های ابتدایی و انتهایی کلمه خطا می‌کند، که ممکن است فاصله و یا علامت‌های نقطه-گذاری باشند.

در نتیجه آزمایش‌های مختلف Ludovik and Zacharski، عملکرد این الگوریتم برای تشخیص زبان متون تک-زبانه شامل زبان‌های عربی و فارسی وقتی خوب است که اندازه متن تحت بررسی بین 500 تا 1000 بایت باشد. نرخ خطای آن 0.27% برای متون 1000 بایتی، 0.52% برای متون 500 بایتی، 2.02% برای متون 100 بایتی، 4.01% برای متون 50 بایتی و 11.92% برای متون 20 بایتی گزارش شده است. باید توجه داشت که این الگوریتم تنها یک ششم نرخ خطای الگوریتم [6] را دارد.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیر پروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	

6. نتیجه گیری

در پژوهش حاضر مشکل تشخیص هوشمند جمله فارسی را در پیکره های مخلوط فارسی و عربی به عنوان یک مسئله تشخیص زبان در حوزه رده بندی متن که از فنون یادگیری ماشین استفاده می کنند نگریسته شد. آنگاه با مروری بر نمونه ای از این پیکره ها، یعنی برخی متون تفسیر قرآن کریم، مشخص شد با وجود اینکه چنین پیکره هایی حجم عظیمی را به خود اختصاص داده اند و در داخل بخش هایی که بطور عمدۀ فارسی هستند تنها بخش های کوچکی با متون عربی مخلوط شده اند، اما این بخش های عربی همیشه به صورت جملات مستقل ظاهر نشده بلکه در قالب عبارت و یا حتی کلمه تنها نیز به کار رفته اند. این امر الگوریتم های قبلی را که با تکیه بر خود واژه ها و آن هم در قالب جملات از قبل از هم تفکیک شده قادر به کار بودند را به چالش کشید. در نهایت الگوریتم تشخیص زبان Ludovik and Zacharski که برای متون چند زبانه طراحی شده است به عنوان چاره کار معرفی شد. این الگوریتم، با تکیه بر گزینشی خاص از ترکیبی خاص از واژه ها و با استفاده از ترکیبی خاص از فنون یادگیری ماشین با ناظر، پس از تقطیع متون چند زبانه به متون تک زبانه می تواند زبان آنها را در بخش هایی با طول تنها دو کلمه متوسط الی سه کلمه کوتاه تشخیص دهد. الگوریتم فوق که روی متون مخلوط فارسی و عربی نیز امتحان شده است هم در تقطیع متون چند-زبانه به متون تک-زبانه و هم در تشخیص زبان این قطعه ها بسیار موفق تر از الگوریتم های قبلی گزارش شده است. آنچه حتی در این الگوریتم هنوز حل نشده است تشخیص زبان بخش های تک-کلمه ای است که با حجمی قابل توجه در پیکره های مخلوط فارسی و عربی به کار رفته اند. امید است در پژوهش های آتی بتوان الگوریتمی طراحی کرد که بتواند الگوریتم Ludovik and Zacharski را تکمیل نماید.

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیرپروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	

7. پایگاه های on-line برای تشخیص زبان متن

در صورتی که استفاده از وبگاههای اینترنتی برای کاربر مشکل ساز نباشد و بتواند به نرخ خطای متوسطی بسنده نماید، نشانیهای اینترنتی ذیل وبگاههایی هستند که وظیفه تشخیص زبان متن را، به طور عمده در متونی که تک-زبانه شدهاند، انجام می‌دهند.

<http://rali.iro.umontreal.ca/Silc/index.jsp?lang=fr>



<http://www.eidetica.com/services/guesser>

<http://www.fuzzums.nl/~joost/talenknobbel/>

<http://odur.let.rug.nl/~vannoord/TextCat/Demo/textcat.html>

<http://www.xrce.xerox.com/competencies/content-analysis/tools/guesser.en.html>

<http://labs.translated.net/language-identifier/>

	عنوان پروژه: فاز اول طرح جامع پیکره زبان فارسی با موضوع فاز اول مطالعاتی ایجاد پیکره متنی زبان فارسی		
	عنوان زیرپروژه: مطالعه و بررسی روشهای تشخیص هوشمند جمله فارسی از عربی در پیکره های مخلوط فارسی و عربی		
	تاریخ: 1388/04/19	ویرایش: 1/0	

مراجع

- [1] نرم افزار "جامع التفاسیر". مرکز تحقیقات کامپیوتری علوم اسلامی، نور. 1385.
- [2] George Forman. Feature Selection for Text Classification. Computational Methods of feature Selection. CRC Press/Taylor and Francis Group. 2007.
- [3] Lena Grothe, Ernesto William De Luca and Andreas Nürnberger. A Comparative Study on Language Identification Methods. In *Proceedings of LREC '08*. 2008.
- [4] G. Grefenstette. Comparing Two Language Identification Schemes. In *Proceedings of 3rd International Conference on Statistical Analysis of Textual Data*. 1994.
- [5] M.J. Martino and R.C. Paulsen. Natural Language Detemination Using Partial Words. In *U.S. Patent No. 6216102 B1*.
- [6] W.B. Cavner and J.M. Trenkle. N-gram Based Text Categorization. In *Proceedings of the Third Annual Symposium on Document Analysis and Information Retrieval*. 1994.
- [7] Yevgeny Ludovik and Ron Zacharski. Multilingual Document Language Recognition for Creating Corpora. In *Proceedings of MT Summit VII*. 1999.
- [8] Chew Lim Tan, Peck Yoke Leong and Shoujie He. Language Identification in Multilingual Documents. In *Proceedings of International Symposium on Intelligent Multimedia*. 1999.
- [9] Yo Ehara and Kumiko Tanaka-Ishii. Multilingual Text Entry Using Automatic Language Detection. In *Proceedings of IJCNLP*. 2008.
- [10] محمود شکرالهی فر، مهدی عنایتی، اردلان حسن زاده، بهاره نظرزاده. سامانه تشخیص زبان متون فارسی و عربی. پژوهش منتشر نشده در دانشگاه نبی اکرم (ص). 1388.