

PersianFarm: دیتاستی برای تطبیق آنتولوژی‌های فارسی

هادی تابع الحجه^۱، بیتا شادگار^۲

^۱ دانشجوی کارشناسی ارشد، گروه مهندسی کامپیوتر، دانشکده مهندسی، دانشگاه شهیدچمران اهواز، اهواز،

h-tabealhojeh@mscstu.scu.ac.ir

^۲ استادیار، گروه مهندسی کامپیوتر، دانشکده مهندسی، دانشگاه شهیدچمران اهواز، اهواز،

Bita.shdgar@scu.ac.ir

چکیده

این مقاله، برای اولین بار مسئله تطبیق آنتولوژی فارسی را به صورت عملی بررسی می‌کند. تطبیق آنتولوژی در پیاده‌سازی وب‌معنایی نقشی کلیدی دارد. با وجود تلاش پژوهشگران برای ساخت آنتولوژی‌های فارسی، متاسفانه تلاشی برای طراحی تطبیق‌گرهای آنتولوژی فارسی صورت نگرفته است. شاید بتوان بزرگترین چالش در طراحی تطبیق‌گرهای فارسی را عدم وجود مجموعه داده‌ی محک^۱ (دیتاست) فارسی دانست. در این مقاله یک مجموعه داده محک استاندارد و جامع برای تطبیق فارسی-فارسی معرفی می‌شود. مجموعه داده محک فارسی شامل هفت آنتولوژی فارسی و یازده جفت تطبیق بین آنتولوژی‌ها است و مطابق با مجموعه داده محک^۲ OntoFarm مربوط به کمپین^۳ OAEI، و براساس استانداردهای آن ایجاد شده است. در ادامه مقاله، عملکرد تعدادی از مهم‌ترین معیارهای محاسبه شباهت رشته‌ای (تطبیق‌گرهای رشته‌ای)، در تطبیق آنتولوژی‌های فارسی بررسی و ارزیابی می‌شوند. نتایج ارزیابی‌ها به صورت معیارهای ارزیابی دقت، فراخوانی و معیار F ارائه شده است.

کلمات کلیدی

تطبیق آنتولوژی فارسی، هم‌ترازی آنتولوژی فارسی، مجموعه داده محک تطبیق آنتولوژی فارسی، معیارهای شباهت معنایی

تطبیق مطلوب باید توانایی رویارویی با انواع ناهمگونی‌ها در منابع آنتولوژی را داشته باشد.

۱- مقدمه

در سال‌های اخیر تعداد زیادی از آنتولوژی‌ها در زبان‌های طبیعی دیگر طراحی و به کار گرفته شده‌اند. این آنتولوژی‌ها ارتباط بسیار مهمی بین اطلاعات موجود در وب‌معنایی و کاربران که اکثراً ترجیح می‌دهند اطلاعات خود را به زبان محلی وارد کنند، برقرار می‌کنند [1]. ظهور آنتولوژی‌ها با زبان‌های مختلف، تطبیق آنتولوژی را وارد مرحله جدیدی کرده است.

در این میان زبان فارسی به‌عنوان زبان رسمی سه کشور ایران، افغانستان و تاجیکستان، از این قاعده مستثنی نیست. اگرچه در سال‌های اخیر استفاده از آنتولوژی‌ها به زبان فارسی روبه‌فزونی است اما عدم وجود سیستم تطبیق مناسب، یکی از موانع پیش‌روی نرم‌افزارهای مبتنی بر آنتولوژی فارسی است. متاسفانه تطبیق آنتولوژی‌های فارسی با چالش‌های متعددی روبه‌رو است. یکی از چالش‌های مهم در طراحی تطبیق‌گرهای فارسی-فارسی، عدم وجود

امروزه آنتولوژی‌ها در علوم مختلفی هم‌چون هوش مصنوعی، وب‌معنایی، مهندسی نرم‌افزار، انفورماتیک پزشکی، علوم کتابداری و پردازش زبان طبیعی استفاده می‌شوند. آنتولوژی به‌عنوان پایگاه دانش، پیچیدگی را کاهش و اطلاعات را سازماندهی می‌کند. در حقیقت، آنتولوژی‌ها توصیفی از پدیده‌های موجود در جهان و روابط بین آن‌ها را فراهم می‌کنند به طوری که برای ماشین قابل تفسیر باشند. در نتیجه ماشین می‌تواند به استنتاج‌های پیچیده‌تری دست یافته و بسیاری از کارهای انسان را به‌عهده بگیرد. ایجاد آنتولوژی‌ها توسط افراد و سازمان‌های متفاوت، ناهمگونی بین آنتولوژی‌ها را در پی خواهد داشت. یعنی یک دامنه دانش توسط چند آنتولوژی بیان می‌شود که هر کدام از اصطلاحات متفاوتی برای تعریف یک موجودیت استفاده می‌کنند. به‌همین دلیل تطبیق آنتولوژی نقشی کلیدی در موفقیت وب‌معنایی دارد. یک سیستم

مجموعه داده محک استاندارد برای ارزیابی کارایی تطبیق‌گرها و مقایسه نتایج آن‌ها با یکدیگر است.

رویکردهای متفاوتی برای ایجاد مجموعه داده محک، توسط پژوهشگران دنبال شده است. ترجان و همکارانش [2] و همچنین فو و همکاران [3] از ترجمه داده‌های محک انگلیسی برای ایجاد مجموعه داده محک غیرانگلیسی استفاده کرده‌اند. عیب اصلی این روش در این است که در تطبیق بین زبانی، دو آنتولوژی ساختاری یکسان خواهند داشت و تطبیق‌گرهای ساختاری به‌صورت دقیق قابل ارزیابی نیستند. رویکرد دیگر، ایجاد داده‌های محک بر پایه آنتولوژی‌های چندزبانی است. این رویکرد کمی پیچیده‌تر است. در این روش از آنتولوژی‌هایی که هر کدام با زبان‌های طبیعی مختلف، حوزه یکسانی (مثلا پزشکی) را توصیف می‌کنند، استفاده می‌شود. در این رویکرد تطبیق‌های مرجع موجود نیستند. یافتن تطبیق‌های مرجع نیازمند استخدام چندین متخصص و صرف زمان زیادی است [4]. یانگ و همکارانش برای یافتن تطبیق‌های مرجع، از یک آنتولوژی میانی و تطبیق‌های مرجع آن استفاده کرده‌اند [5]. به‌عنوان نمونه تطبیق بین زبان‌های پرتغالی و اسپانیایی را می‌توان با تطبیق بین دو آنتولوژی از زبان‌های مذکور و آنتولوژی زبان سوم مثلا انگلیسی را به‌دست آورد. رویکرد دیگر برای ایجاد مجموعه داده محک، استفاده از منابعی غیر از آنتولوژی است که به چند زبان در دسترس هستند. است. مثلا می‌توان از دایرکتوری‌های یاهو که در زبان‌های مختلف موجود است، استفاده کرد. برای ساخت مجموعه محک، ابتدا دایرکتوری‌ها به آنتولوژی تبدیل می‌شوند. عبارت‌های زبان‌های مختلف که یک دایرکتوری را بیان می‌کنند، منطبق هستند و تطبیق‌های مرجع محسوب می‌شوند. این رویکرد در [6] استفاده شده است. متاسفانه آنتولوژی‌های ساخته شده در این روش دارای ساختار بسیار ساده هستند و مفاهیم آنتولوژی ارتباطات کم و ناقصی دارند.

کمپین OAEI در سال ۲۰۱۲ میلادی، یک مجموعه محک به نام MultiFarm برای تطبیق بین‌زبانی آنتولوژی معرفی کرد. این مجموعه محک، از ترجمه مفاهیم آنتولوژی‌های مجموعه محک OntoFarm به نه زبان از جمله زبان‌های روسی، فرانسوی، چینی و پرتغالی ایجاد شده است [4].

چالش مهم دیگر، انتخاب معیارهای شباهت و ترکیب مناسب آن‌ها است. در زمینه انتخاب معیارهای شباهت، بیشتر پژوهش‌های انجام شده مربوط به تطبیق‌گرهای انگلیسی زبان است. چیتهم و همکارانش تعدادی از معیارهای رشته‌ای را ارزیابی کرده‌اند [7]. آن‌جی‌او و همکارانش نیز برای طراحی تطبیق‌گر خود انواع معیارهای مشابهت معنایی را مورد ارزیابی قرار داده است [8] و در گام بعد، معیارهایی که بهترین نتایج را داشته‌اند را ترکیب کرده است.

این مقاله موضوع مذکور را بحث و بررسی می‌کند. بخش دوم، نحوه ایجاد مجموعه داده محک (دیناست) فارسی را تشریح می‌کند. در بخش سوم، عملکرد معیارهای مشابهت رشته‌ای بر اساس آنتولوژی‌های فارسی ارزیابی و مقایسه می‌شود. در پایان بخش چهارم به نتیجه‌گیری می‌پردازد.

۲- مجموعه داده محک فارسی

از آنجا که کارایی متدهای تطبیق به ذات آنتولوژی‌ها وابسته است، ارزیابی اصولی و سیستماتیک متدهای تطبیق برای روشن کردن نقات ضعف و قوت

سیستم‌های تطبیق‌گر، ضروری است، از سال ۲۰۰۴ میلادی، کمپین OAEI ارزیابی سیستم‌های تطبیق آنتولوژی را انجام می‌دهد. کمپین OAEI چندین مجموعه داده محک را فراهم کرده است که هر کدام عمل تطبیق را از جنبه خاصی ارزیابی می‌کند. هر مجموعه داده شامل چند آنتولوژی و تطبیق‌های مرجع بین آن‌ها است.

در این پژوهش یک مجموعه داده برای تطبیق فارسی-فارسی معرفی می‌شود. مجموعه داده محک پیشنهادی، با رویکرد ترجمه مجموعه داده‌های محک انگلیسی و از ترجمه یکی از مجموعه داده‌های انگلیسی کمپین OAEI به نام OntoFarm ایجاد شده است. که شامل هفت آنتولوژی مختلف است و تطبیق‌های متقابل بین آنتولوژی‌ها به‌صورت دستی توسط متخصصان مشخص شده‌اند. آنتولوژی‌ها به زبان فارسی ترجمه شده‌اند. نحوه ترجمه آنتولوژی انگلیسی زبان به زبان فارسی دارای چالش‌هایی است که در فرایند تطبیق تأثیرگذار خواهند بود. در ادامه ویژگی‌های مجموعه فارسی شرح داده می‌شود.

۱. **نام موجودیت‌های آنتولوژی:** نام یک موجودیت (کلاس، نمونه، ویژگی و...) می‌تواند بخشی از شناسه URI باشد مانند: موضوع http://cmt_per# و یا به‌صورت برچسب و به‌عنوان متادیتا تعریف شود مانند:

<rdfs:label xml:lang="per"> موضوع </rdfs:label>
قراردادن نام در شناسه URI، باعث می‌شود طراح به‌راحتی موجودیت‌ها را به خاطر بسپارد. از طرفی شناسه باید مختصر باشد و برای بسیاری از موجودیت‌ها مناسب نیست. برای طراحی مجموعه داده محک فارسی، نام‌ها درون برچسب تعریف شده‌اند. همچنین زبان عبارت‌ها به‌وسیله تگ زبان ("xml:lang="per") برچسب خورده‌اند.

۲. **کدگذاری آنتولوژی‌ها:** در آنتولوژی‌های ترجمه شده از کدگذاری UTF-8 استفاده شده است. اهمیت تعیین نوع کدگذاری در آن است که ممکن است بعضی از سیستم‌های تطبیق‌گر فقط با نوع خاصی از کدگذاری‌ها سازگار باشد. و این پارامتر باید مدنظر برنامه‌نویس قرار گیرد.

۳. **جداسازی کلمات:** ممکن است نام یک موجودیت، عبارتی چندکلمه‌ای باشد. معمولا کلمات با کاراکترهای مخصوصی نظیر کاراکترهای فاصله (space)، کاما، خط‌تیره (dash) و یا خط‌زیر (underline) از هم جدا می‌شوند. کاراکتر خط‌تیره برای جداسازی کلمات در آنتولوژی‌های فارسی انتخاب شده است. جداسازی به مقایسه موثر نام موجودیت‌ها در تطبیق کمک می‌کند. مزیت دیگر جداسازی کلمات یک عبارت، تشخیص نام‌های چندبخشی است. مثلا «نرم‌افزار» نامی دو بخشی است و میان دو بخش آن، خط‌تیره قرار نمی‌گیرد.

۴. **ترجمه مفاهیم:** مهم‌ترین قسمت ایجاد آنتولوژی‌های محک، ترجمه دقیق و مناسب است. روند ترجمه عبارات در ادامه شرح داده می‌شود. ابتدا کلمات اختصاری بسط داده شده‌اند. به‌عنوان مثال موجودیت «PC» به «Program Chair» بسط داده شده و سپس ترجمه شده است. ترجمه عبارات با توجه به مفهوم موجودیت‌های آنتولوژی انجام شده است و از ترجمه تحت‌لفظی اجتناب شده است. مثلا موجودیت «camera ready paper» به «نسخه نهایی مقاله» ترجمه شده است و نه «مقاله دوربین آماده». همچنین تمامی اسم‌ها به‌صورت مفرد قرار

۳- مقایسه تطبیق گره‌های رشته‌ای در تطبیق آنتولوژی‌های فارسی

معمولا تطبیق گره‌ها از مجموعه‌ای از معیارهای مشابهت شامل معیارهای رشته‌ای، ساختاری و زبانی برای مقایسه موجودیت‌های دو آنتولوژی استفاده می‌کنند. اکثر سیستم‌های تطبیق طراحی شده، برای کاهش محاسبات از معماری سری برای ترکیب انواع تطبیق گره‌ها استفاده می‌کنند. در این معماری، ابتدا تمام موجودیت‌های آنتولوژی اول با تمام موجودیت‌های آنتولوژی دوم به‌وسیله معیارهای رشته‌ای مقایسه می‌شوند. سپس موجودیت‌هایی که تشابه آن‌ها از یک مقدار آستانه بیشتر باشد، با معیارهای زبانی و ساختاری سنجیده می‌شوند. بنابراین معیارهای رشته‌ای در مقایسه نام موجودیت‌ها از اهمیت خاصی برخوردارند. همچنین مقدار آستانه، پارامتری تاثیرگذار در فرآیند تطبیق است.

در این پژوهش بیش از ۲۰ معیار رشته‌ای از رابط‌های برنامه‌نویسی Secondstring^۷ و Simmetrics^۷ ارزیابی می‌شوند. هر معیار ۱۱ مرتبه برای یافتن تطبیق‌های ۱۱ جفت آنتولوژی فارسی (مجموعه داده محک، ۱۱ جفت تطبیق مرجع دارد) آزمایش شده است. روند پیاده‌سازی فرآیند تطبیق در ادامه شرح داده می‌شود. ابتدا موجودیت‌های دو آنتولوژی به‌وسیله رابط برنامه‌نویسی Alignment API استخراج می‌شوند. سپس شباهت موجودیت‌ها به‌صورت دوجه‌دو محاسبه می‌شود. اندازه‌ی شباهت حاصل، مقداری نرمال در بازه‌ی صفر و یک است. برای ارزیابی نتایج از معیارهای دقت^۸، فراخوانی^۹ و معیار F^{۱۰} استفاده شده است.

$$\text{Precision} = \frac{|\text{alignment given} \cap \text{correct alignment}|}{|\text{alignment given}|} \quad (1)$$

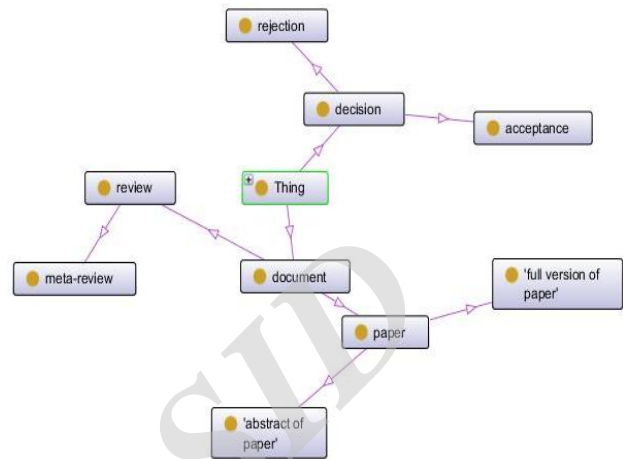
$$\text{Recall} = \frac{|\text{alignment given} \cap \text{correct alignment}|}{|\text{correct alignment}|} \quad (2)$$

$$\text{Fmeasure} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (3)$$

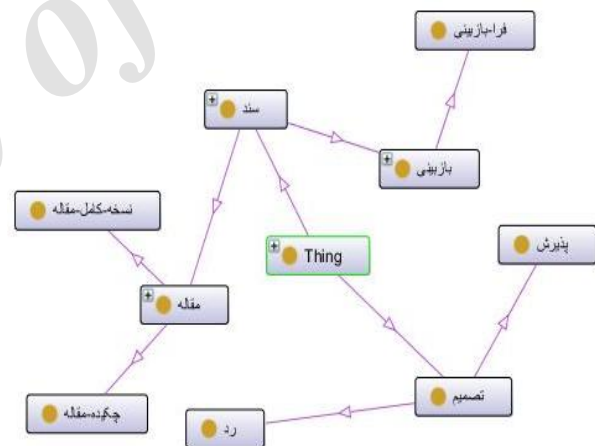
معیارهای ارزیابی به‌ازای مقادیر آستانه مختلف محاسبه شده‌اند. شکل (۳) نتایج ارزیابی معیارهایی که کارایی بهتری داشتند را نمایش می‌دهد. میانگین ارزیابی‌های هر معیار شباهت، به‌ازای ۱۱ فرآیند تطبیق، به‌ازای آستانه‌های متفاوت محاسبه شده است. شکل (۴) نیز میانگین ماکزیمم معیار F به‌ازای هر تطبیق‌گر را نمایش می‌دهد. برای پیاده‌سازی فرآیند تطبیق از زبان برنامه‌نویسی جاوا استفاده شده است.

معیار لوشترین با معیار F برابر با ۶۶ درصد بهترین نتیجه را داشته است. همچنین معیار لوشترین که بر مبنای فاصله ویرایشی عمل می‌کند، در مقایسه با سایر معیارهای مبتنی بر فاصله ویرایشی مانند اسمیت-واترمن و نیدلمن وانچ، بهتر عمل کرده است زیرا معیار لوشترین به هر سه عمل ویرایشی (حذف، جایگذاری و تعویض) وزن برابری نسبت می‌دهد و به پیش‌وند یا پس-وند مشترک حساسیت ندارد. به‌عنوان نمونه لوشترین به برچسب‌های «مقاله» و «فراخوان-مقاله» به‌درستی امتیاز مشابهت کمتری می‌دهد. زیرا این دو موجودیت اگرچه در عبارت «مقاله» مشترک‌اند، اما تطابق ندارند. در میان معیارهای مبتنی بر توکن‌بندی رشته‌ها، qgram در مقایسه با سایر متدها مانند

داده شده و حرف «ی» اضافه حذف شده است. مثلا عبارت «نسخه‌ی کامل مقاله‌ها» به «نسخه-کامل-مقاله» تبدیل شده است. شکل (۱) بخشی از آنتولوژی cmt از مجموعه داده انگلیسی کنفرانس را نمایش می‌دهد و شکل (۲) آنتولوژی ترجمه شده آن را نمایش می‌دهد. جدول ۱ نیز تعداد موجودیت‌های هر آنتولوژی را نشان می‌دهد. مجموعه داده فارسی شامل ۷ آنتولوژی و ۱۱ جفت تطبیق مرجع است.



شکل (۱): بخشی از آنتولوژی cmt-en



شکل (۲): بخشی آنتولوژی cmt-per که با شکل (۲) متناظر است

جدول (۱): تعداد موجودیت‌های آنتولوژی‌های مجموعه داده محک فارسی

نام آنتولوژی	تعداد		
	تعداد کلاس	ویژگی‌های نوع داده	تعداد ویژگی شیء
Conference	۶۰	۱۸	۴۶
Cmt	۳۶	۳۶	۴۹
Ekaw	۷۴	۰	۳۳
Idas	۱۰۴	۲۰	۳۰
Iasted	۱۴۰	۳	۳۸
Confof	۳۸	۲۳	۱۳
Sigkdd	۴۹	۱۱	۱۷

بین آنها است. در بخش دوم مقاله، معیارهای مشابهت رشته‌ای به‌وسیله دیتاست معرفی شده مورد ارزیابی قرار گرفتند. نتایج حاصل نشان می‌دهد که معیار لونشتین با مقدار ۶۶ درصد معیار F بهترین نتیجه را داشته است.

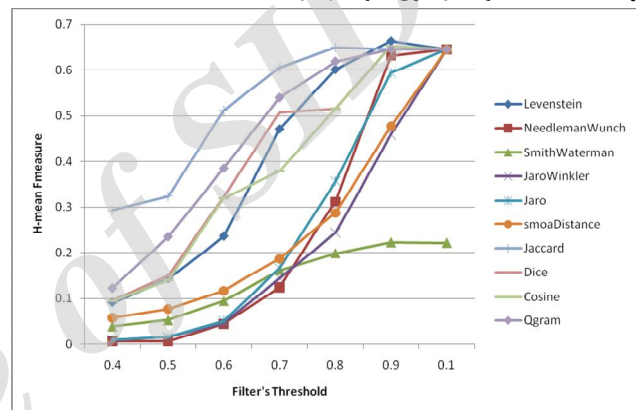
مراجع

- [1] Fu, Bo, Rob Brennan, and Declan O'sullivan. "Cross-Lingual Ontology Mapping and Its Use on the Multilingual Semantic Web." *MSW*. 2010.
- [2] Trojahn, Cássia, Paulo Quaresma, and Renata Vieira. "A framework for multilingual ontology mapping." (2008).
- [3] Fu, Bo, Rob Brennan, and Declan O'Sullivan. "Cross-lingual ontology mapping—an investigation of the impact of machine translation." *The Semantic Web*. Springer Berlin Heidelberg, 2009. 1-15.
- [4] Meilicke, Christian, et al. "MultiFarm: A benchmark for multilingual ontology matching." *Web Semantics: Science, Services and Agents on the World Wide Web* 15 (2012): 62-68.
- [5] Jung, Jason J., Anne Håkansson, and Ronald Hartung. "Indirect alignment between multilingual ontologies: A case study of korean and swedish ontologies." *Agent and Multi-Agent Systems: Technologies and Applications*. Springer Berlin Heidelberg, 2009. 233-241.
- [6] Caraciolo, Caterina, et al. "Results of the ontology alignment evaluation initiative 2008." *Proc. 3rd ISWC workshop on ontology matching (OM)*. 2008.
- [7] Cheatham, Michelle, and Pascal Hitzler. "String similarity metrics for ontology alignment." *The Semantic Web—ISWC 2013*. Springer Berlin Heidelberg, 2013. 294-309.
- [8] Ngo, DuyHoa, Zohra Bellahsene, and Konstantin Todorov. "Opening the black box of ontology matching." *The Semantic Web: Semantics and Big Data*. Springer Berlin Heidelberg, 2013. 16-30.

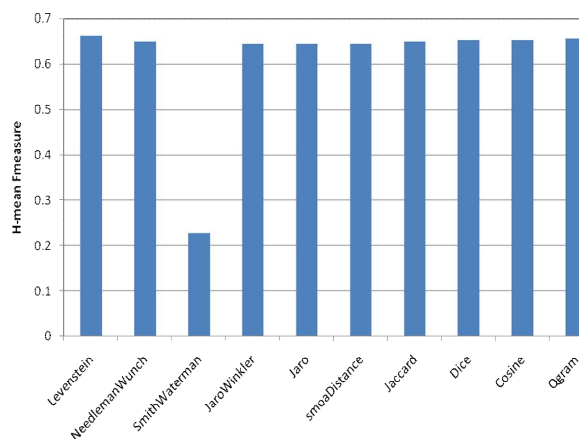
زیر نویس‌ها

- 1 Benchmark dataset
- 2 <http://nb.vse.cz/~svatek/ontofarm.html>
- 3 Ontology Alignment Evaluation Initiative
- 4 Multilingual ontology matching
- 5 Reference Alignment
- 6 <http://secondstring.sourceforge.net>
- 7 <http://sourceforge.net/projects/simmetrics>
- 8 Precision
- 9 Recall
- 10 F-measure

dice و cosine نتیجه بهتری داشته است. روند مهم دیگر این است که اگرچه بیشینه تشخیص اکثر معیارها در شکل (۴)، بسیار به هم نزدیک است اما سه معیار لونشتین، جاکارد و qgram در مقادیر آستانه کمتر از عدد یک به بیشینه معیار F دست پیدا کرده‌اند. در نتایج حاصل دو روند مهم مشاهده می‌شود. اول این که با افزایش مقدار آستانه، فراخوانی کاهش می‌یابد که این امر منجر به افزایش هم‌زمان دقت می‌شود. مقدار آستانه ۱ هیچ تطبیقی را فیلتر نمی‌کند و یک نقطه شکستگی در نمودار خود دارند. به عبارت دیگر با انتخاب این معیارها علاوه بر دستیابی به نتیجه مناسب، بسیاری از جفت موجودیت‌های غیرمنطبق فیلتر می‌شوند. همان‌طور که قبلاً ذکر شد، انتخاب مقدار آستانه از آنجا اهمیت می‌یابد که در سیستم‌های تطبیق متوالی (سری)، در مرحله دوم فقط مشابهت معنایی جفت موجودیت‌هایی بررسی می‌شوند که مقدار مشابهت آنها در مرحله اول از مقدار آستانه بیشتر باشد. این روش در تطبیق آنتولوژی‌های بزرگ کاربرد بسیاری دارد زیرا زمان اجرای کمتری دارد و به حافظه کمتری در طول اجرا نیاز دارد.



شکل (۳): میانگین معیار F مربوط به تطبیق گره‌های رشته‌ای در آستانه‌های متفاوت



شکل (۴): میانگین بیشینه معیار F مربوط به تطبیق گره‌های رشته‌ای

۴- نتیجه

در این مقاله برای اولین بار یک دیتاست فارسی تشکیل شده از تعدادی آنتولوژی فارسی برای طراحی و ارزیابی سیستم‌های تطبیق آنتولوژی فارسی معرفی شده است. دیتاست معرفی شده شامل ۷ آنتولوژی و یازده جفت تطبیق