



دانشکده مهندسی کامپیوتر

# کاربست یک الگوریتم پیونددهی موجودیت برای مسئله‌ی ابهام‌زدایی معنی واژگان در زبان فارسی

پایان‌نامه یا رساله برای دریافت درجه کارشناسی  
در رشته مهندسی کامپیوتر گرایش نرم‌افزار - هوش

نام دانشجو

بنفشه کریمیان

استاد راهنما:

دکتر مینایی

تیر ماه ۱۳۹۸

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

## تأییدیه‌ی هیأت داوران جلسه‌ی دفاع از پایان‌نامه/رساله

نام دانشکده: دانشکده کامپیوتر

نام دانشجو: بنفشه کریمیان

عنوان پایان‌نامه یا رساله: کاربرست یک الگوریتم پیونددهی موجودیت برای مسئله‌ی ابهام‌زدایی معنی‌واژگان

در زبان فارسی

تاریخ دفاع: ۱۳۹۸/۳/۲۶

رشته: مهندسی کامپیوتر

گرایش: نرم‌افزار - هوش

ردیف	سمت	نام و نام خانوادگی	مرتبه دانشگاهی	دانشگاه یا مؤسسه	امضا
۱	استاد راهنما	دکتر بهروز مینائی	دانش‌یار	دانشگاه علم و صنعت	
۲	استاد راهنما				
۳	استاد مشاور				
۴	استاد مشاور				
۵	استاد مدعو خارجی				
۶	استاد مدعو خارجی				
۷	استاد مدعو داخلی	دکتر سید صالح اعتمادی	استادیار	دانشگاه علم و صنعت	
۸	استاد مدعو داخلی				

## تأییدیه‌ی صحت و اصالت نتایج

### باسمه تعالی

اینجانب بنفشه کریمیان به شماره دانشجویی ۹۴۵۲۱۱۸۹ دانشجوی رشته مهندسی کامپیوتر مقطع تحصیلی کارشناسی تأیید می‌نمایم که کلیه‌ی نتایج این پایان‌نامه/رساله حاصل کار اینجانب و بدون هرگونه دخل و تصرف است و موارد نسخه‌برداری شده از آثار دیگران را با ذکر کامل مشخصات منبع ذکر کرده‌ام. در صورت اثبات خلاف مندرجات فوق، به تشخیص دانشگاه مطابق با ضوابط و مقررات حاکم (قانون حمایت از حقوق مؤلفان و مصنفان و قانون ترجمه و تکثیر کتب و نشریات و آثار صوتی، ضوابط و مقررات آموزشی، پژوهشی و انضباطی ...) با اینجانب رفتار خواهد شد و حق هرگونه اعتراض درخصوص احقاق حقوق مکتسب و تشخیص و تعیین تخلف و مجازات را از خویش سلب می‌نمایم. در ضمن، مسئولیت هرگونه پاسخگویی به اشخاص اعم از حقیقی و حقوقی و مراجع ذی‌صلاح (اعم از اداری و قضایی) به عهده‌ی اینجانب خواهد بود و دانشگاه هیچ‌گونه مسئولیتی در این خصوص نخواهد داشت.

نام و نام خانوادگی: بنفشه کریمیان

امضا و تاریخ: ۱۳۹۸/۳/۲۶

## مجوز بهره‌برداری از پایان‌نامه

بهره‌برداری از این پایان‌نامه در چهارچوب مقررات کتابخانه و با توجه به محدودیتی که توسط استاد راهنما به شرح زیر تعیین می‌شود، بلامانع است:

☐ بهره‌برداری از این پایان‌نامه/ رساله برای همگان بلامانع است.

✓ بهره‌برداری از این پایان‌نامه/ رساله با اخذ مجوز از استاد راهنما، بلامانع است.

☐ بهره‌برداری از این پایان‌نامه/ رساله تا تاریخ ..... ممنوع است.

نام استاد یا اساتید راهنما:

دکتر بهروز مینائی

تاریخ:

امضا:

## چکیده

ابهام‌زدایی معنای واژگان یکی از مهم‌ترین ارکان استخراج دانش از متن به‌شمار می‌آید، که در آن هر واژه‌ی یک متن به یک هم‌نشیم از یک وردنت پیوند داده می‌شود. این مسئله بسیار مشابه با مسئله‌ی پیونددهی موجودیت‌ها می‌باشد که در آن هر کلمه به یک موجودیت در گراف دانش یا پایگاه دانش پیونددهی می‌شود. در فارسی این دو مسئله با استفاده از وردنت فارسی (فارس‌نت) و گراف دانش فارسی (فارس‌بیس) پوشش داده شده‌اند. با این وجود، در حالی که یک رویکرد متحد برای ابهام‌زدایی معنای واژگان و پیونددهی موجودیت‌ها برای زبان انگلیسی ارائه داده شده است و گراف‌دانش انگلیسی به وردنت انگلیسی نگاشت شده است، در زبان فارسی این رویکرد متحد و نگاشت اتخاذ نشده است.

این پروژه بر آن است که نگاشت بین فارس‌نت و فارس‌بیس که خلأ آن در زبان فارسی موجود است را پوشش دهد. در این راستا، ابتدا بر اساس یک الگوریتم پیونددهی موجودیت (وابسته به فارس‌بیس) یک الگوریتم مشابه برای ابهام‌زدایی معنایی (بر اساس فارس‌نت) ارائه می‌شود. در ادامه، برای گسترش ابهام‌زدایی واژگان در زبان فارسی با الهام‌گیری از پروژه‌ی ببیل‌نت<sup>۱</sup> به ترکیب توأمان فارس‌نت و فارس‌بیس می‌پردازیم؛ به این صورت که الگوریتم ابهام‌زدای ارائه شده را به طور توأمان با الگوریتم پیونددهی موجودیت مذکور بر یک مجموعه داده اعمال می‌نماییم، تا هر کدام از واژگان مجموعه داده به هم‌نشیم‌های فارس‌نت و موجودیت‌های فارس‌بیس پیوند داده شوند. سپس هم‌نشیم‌ها و موجودیت‌های متناظر به یک‌دیگر نگاشت شده و فراوانی این نگاشت اندازه‌گیری می‌شود. بر اساس این فراوانی، و با قراردادن یک آستانه<sup>۲</sup> هم‌نشیم‌هایی که تعداد باهم‌آیی آن‌ها از مقدار آستانه بیشتر است را با هم ادغام می‌نماییم. در نهایت، با مقایسه با یک دادگان طلایی<sup>۳</sup> معیارهای دقت و بازخوانی خروجی ارائه شده را محاسبه می‌نماییم. شایان ذکر است که دادگان طلایی نگاشت شده توسط یک شخص خبره تولید شده و شامل ۷۶۷ نگاشت است. معیار دقت، بازخوانی و تعداد نگاشت نهایی به‌دست آمده از این پروژه به ترتیب برابر با ۷۰ درصد، ۳۰ درصد و ۵۰۰۰ هم‌نشیم است.

**واژه‌های کلیدی:** ابهام‌زدایی معنایی، پیونددهی موجودیت، فارس‌نت، فارس‌بیس، ببیل‌نت

---

<sup>1</sup> Babelnet

<sup>2</sup> Threshold

<sup>3</sup> Gold data

## فهرست مطالب

۱	فصل ۱: مقدمه
۲	۱-۱ - مقدمه
۴	فصل ۲: مروری بر منابع
۵	۱-۲ - مقدمه
۵	۲-۲ - تعاریف، اصول و مبانی نظری
۵	۱-۲-۲ - فارسی‌نت
۶	۲-۲-۲ - فارسی‌بیس
۶	۳-۲ - مروری بر ادبیات موضوع
۶	۱-۳-۲ - ابهام‌زدایی معنایی
۹	۲-۳-۲ - پیونددهی موجودیت
۱۰	۳-۳-۲ - نگاشات میان ابهام‌زدایی معنایی و پیونددهی موجودیت
۱۰	۴-۲ - نتیجه‌گیری
۱۱	فصل ۳: روش تحقیق
۱۲	۱-۳ - مقدمه
۱۲	۲-۳ - تشریح کامل روش تحقیق
۱۴	فصل ۴: نتایج و تفسیر آنها
۱۵	۱-۴ - مقدمه
۱۵	۲-۴ - محتوا
۱۵	۱-۲-۴ - مجموعه داده‌ها
۱۵	۲-۲-۴ - ارزیابی نتایج
۱۸	فصل ۵: جمع‌بندی و پیشنهادها
۱۹	۱-۵ - مقدمه
۱۹	۲-۵ - جمع‌بندی
۱۹	۳-۵ - نوآوری
۱۹	۴-۵ - پیشنهادها
۲۱	مراجع
۲۵	فهرست واژگان فارسی به انگلیسی
۲۶	فهرست واژگان انگلیسی به فارسی





## فهرست اشکال

- شکل (۱-۴) مقدار معیار دقت و بازخوانی برای آستانه‌های متفاوت ..... ۱۶
- شکل (۲-۴) تعداد اجزای سخن در میان نگاشت‌ها برای هر آستانه ..... ۱۶
- شکل (۳-۴) تعداد نگاشت‌ها در هر آستانه ..... ۱۷

# فصل ١ :

## مقدمه

## ۱-۱- مقدمه

زبان وسیله‌ی انتقال اطلاعات بین موجودات و نسل‌های مختلف است و این انتقال اطلاعات را می‌توان به دسته‌ی زبانی و غیر زبانی تقسیم کرد. پردازش زبان طبیعی<sup>۱</sup> یک روش برای تحلیل و بررسی زبان متن و یکی از حوزه‌های فعال در هوش مصنوعی است [1]. در این میان فهم معنای یک متن از اهمیت بسزایی برخوردار است و برای مدت زمانی طولانی از اهداف مهم تحقیق‌های این حوزه و حوزه‌های مشابه بوده‌است [2]. دو مسئله‌ی مهم برای رسیدن به فهم معنای یک متن ابهام‌زدایی معنایی<sup>۲</sup> و پیونددهی موجودیت<sup>۳</sup>‌اند. این دو مسئله از مسائل شناخته‌شده در این حوزه‌اند [3].

هر واژه ممکن است معانی متفاوتی با توجه به محتوای مربوطه داشته باشد، که به آن پولیسمی<sup>۴</sup> می‌گویند [4]. برای مثال واژه‌ی شیر در سه جمله‌ی زیر معانی متفاوتی دارد:

- شیر سلطان جنگل است.
- شیر آب باز بود.
- شیر در لیوان ریخته‌شده بود.

ابهام‌زدایی معنایی برای تشخیص معنای یک واژه در یک متن با توجه به محتوای متن از میان معانی مختلف آن واژه است و از مسائل سخت در حوزه‌ی هوش مصنوعی محسوب می‌شود. روش‌های بسیاری برای ابهام‌زدایی معنایی ارائه شده‌اند ولی با این حال این مسئله هنوز حل نشده محسوب می‌شود. با اینکه لغت‌نامه‌ها می‌توانند برای ابهام‌زدایی معنایی مناسب باشند، استفاده از وردنت پیشنهاد شده است [4].

پیونددهی موجودیت نام‌دار همان نگاشت کردن موجودیت‌های موجود در یک متن به موجودیت‌های یک پایگاه‌دانش است [5]. یک موجودیت نام‌دار ممکن است به چند موجودیت در پایگاه‌دانش با یک نام نگاشت شود، که موجب ابهام در پیونددهی می‌شود [6]. برای مثال، شیشه در ویکی‌پدیای فارسی، که نوعی پایگاه‌دانش است، به موجودیت‌های مختلفی نظیر نام فیلم، موادمخدر، یک رمان، به معنای خود شیشه و غیره و در ویکی‌پدیای انگلیسی به خود شیشه، عینک، نام یک روستا، نام یک رود، نام چند فیلم، رمان و غیره نگاشت می‌شود. یک الگوریتم پیونددهی موجودیت می‌بایست با استفاده از محتوای متن و اطلاعات

<sup>1</sup> Natural Language Processing (NLP)

<sup>2</sup> Word Sense Disambiguation (WSD)

<sup>3</sup> Named-Entity Linking

<sup>4</sup> Polysemy

موجود در پایگاه‌دانش بتواند این ابهام را حل کند. رفع این ابهام، در صورتی که یک موجودیت از نام کامل خود استفاده نکند، بسیار مشکل می‌شود [6].

در پیونددهی موجودیت برخلاف ابهام‌زدایی معنایی یک نگاشت ممکن است نیمه باشد، به این معنا که نام کامل یک موجودیت نیآورده شده باشد ولی با استفاده از محتوای متن بتوان آن را نگاشت کرد. با وجود اختلافات موجود، ابهام‌زدایی معنایی و پیونددهی موجودیت درهم‌تنیده‌اند، چرا که هر دو به ابهام‌زدایی بخش‌های یک متن می‌پردازند [2]. با اینکه در زبان انگلیسی این دو به هم پیوند داده شده‌اند، در زبان فارسی تا کنون مجزا بررسی شده‌اند.

در این پروژه روشی برای نگاشت وردنت فارسی یا همان فارس‌نت به گراف دانش موجودیت‌های فارسی یا همان فارس‌بیس ارائه می‌دهیم، که از نتایج دو الگوریتم مشابه برای ابهام‌زدایی معنایی و پیونددهی موجودیت برای این نگاشت استفاده می‌کند.

در بخش‌های آتی، ابتدا به بررسی کارهای پیشین در این حوزه و سپس، به ارائه‌ی دقیق روش ارائه‌شده می‌پردازیم. در انتها نتایج را بررسی و به جمع‌بندی و ارائه‌ی پیشنهادات می‌پردازیم.

## فصل ۲ :

### مروری بر منابع

## ۲-۱- مقدمه

در این فصل ابتدا به بررسی فارس‌نت و فارس‌بیس و اصول آن دو و سپس، به بررسی و مطالعه‌ی کارهای مشابه انجام شده در این حوزه می‌پردازیم. این مطالعات را می‌توان به سه دسته‌ی مطالعات بر روی ابهام‌زدایی معنایی، مطالعات بر روی پیونددهی موجودیت و مطالعات بر روی نگاشت بین این دو تقسیم‌بندی نمود.

## ۲-۲- تعاریف، اصول و مبانی نظری

در این بخش ابتدا به بررسی فارس‌نت و اصطلاحات موجود در آن و سپس به بررسی فارس‌بیس می‌پردازیم.

### ۲-۲-۱- فارس‌نت

فارس‌نت [7] یک پایگاه‌داده‌ی لغوی فارسی و در واقع نسخه‌ی فارسی وردنت [8] است. در فارس‌نت کلمات به صورت گروه‌های هم‌معنایی یا هم‌نشیم‌ها<sup>۱</sup> دسته‌بندی شده‌اند. برای هر گروه هم‌نشیم اطلاعاتی نظیر توصیف معنا یا همان گلاس<sup>۲</sup>، اطلاعات نحوی، مثالی از کاربرد معنی و غیره نگهداری می‌شود. هم‌نشیم‌ها می‌توانند با هم روابط زیر را داشته باشند:

- روابط هایپرینیم/ هایپونیم<sup>۳</sup>: برای مثال "رنگ" هایپرینیم برای "قرمز" و "قرمز" هایپونیم برای "رنگ" است.
- روابط مرونیمی<sup>۴</sup>: برای مثال دست و چشم برای انسان مرونیما.
- روابط تضادی<sup>۵</sup>

---

<sup>1</sup> Synsets

<sup>2</sup> Gloss

<sup>3</sup> Hypernymy/ hyponym

<sup>4</sup> Meronymy

<sup>5</sup> Antonymy

همچنین رابطه‌ی خواهری<sup>۱</sup> را می‌توان برای هم‌نشیم‌ها تعریف کرد، که به معنای داشتن یک هایپر نیم یکسان بین دو هم‌نشیم است.

## ۲-۲-۲ - فارس بیس

گراف‌های دانش، گراف‌های عظیمی از موجودیت‌های با ارتباط داخلی‌اند. فارس بیس [9]، [10] یک گراف دانش فارسی است، که از منابع اطلاعاتی نظیر ویکی‌پدیا درست شده است.

## ۲-۳ - مروری بر ادبیات موضوع

در این بخش ابتدا به بررسی مطالعات بر روی ابهام‌زدایی معنایی و سپس به مطالعات بر روی پیونددهی موجودیت می‌پردازیم. در انتها، مطالعات بر روی نگاشت بین این دو را بررسی می‌کنیم.

## ۲-۳-۱ - ابهام‌زدایی معنایی

همان‌گونه که در بخش ۱ بیان کردیم، ابهام‌زدایی معنایی تشخیص معنای درست یک واژه در یک متن با توجه به محتوای متن از میان معانی مختلف آن واژه است. سه روش کلی برای ابهام‌زدایی معنایی وجود دارد که در ادامه به بررسی مطالعات انجام شده بر هر سه روش می‌پردازیم [1].

### • روش مبتنی بر دانش<sup>۲</sup>

در این روش از یک منبع لغوی خارجی مانند لغت‌نامه‌ها استفاده می‌شود که به دلیل نبود نیاز به یادگیری، مجموعه‌داده‌ی برچسب خورده‌ای لازم نیست. با اینکه روش‌های با ناظر ممکن است بهینه‌تر باشند، این روش از مزیت‌های بسیاری مانند پیاده‌سازی آسان و نبود نیاز به مجموعه‌داده‌ی برچسب دار بهره‌مند است [1]، [11]. اولین روش ارائه شده در این دسته، روش لسک<sup>۳</sup> است که سعی در

---

<sup>1</sup> Sisterhood

<sup>2</sup> Knowledge based

<sup>3</sup> Lesk

بیشینه کردن هم‌پوشانی لغوی بین معانی مختلف کلمات درون محتوا دارد [12]. در این روش معانی کلمات درون یک محتوا به صورت هم‌زمان ابهام‌زدایی می‌شوند. کیلگاریف و روزنویگ [13] روش لسک ساده‌شده<sup>۱</sup> را ارائه دادند، که در هر زمان تنها یک کلمه را ابهام‌زدایی می‌کند. این روش فضای جست‌وجو را کاهش بسزایی داد.

• روش با ناظر<sup>۲</sup>

در روش‌های با ناظر یک طبقه‌بند<sup>۳</sup> با استفاده از یک مجموعه‌داده‌ی برچسب‌زده‌شده، که در اینجا یک کلمه و معنای متناظر آن در یک محتوا است، آموزش داده می‌شود. سپس از مدل یافت شده برای یافتن معنای متناظر کلمات در متن‌های بدون برچسب استفاده می‌شود [11]. در [14] روشی با استفاده از طبقه‌بند بیز ساده<sup>۴</sup> ارائه شده است. لی، نگ و چیا [15] با استفاده از اجزای کلام<sup>۵</sup>، لغات اطراف در محتوا، همایندها<sup>۶</sup> و روابط نحوی<sup>۷</sup> و یک ماشین بردار پشتیبانی<sup>۸</sup> ابهام‌زدایی معنایی را انجام دادند. هم‌چنین در [16] نیز از بیز ساده برای ابهام‌زدایی استفاده شده است. در [17] و [18] به ترتیب از شبکه‌ی عصبی<sup>۹</sup> و لیست تصمیم<sup>۱۰</sup> به عنوان طبقه‌بند برای ابهام‌زدایی استفاده شده است. لیست تصمیم ارائه‌شده برای زبان فرانسه و اسپانیایی است و از ویژگی‌های اجزای کلام و فاصله‌ی کلمات در متن استفاده می‌کند. علاوه بر این مطالعات، چترجی و میسرا [19] یک مدل قابل آموزش ارائه داده‌اند، که با گرفتن مجموعه‌داده‌ی برچسب‌زده‌شده یادگیری را انجام داده و برای یک مجموعه‌داده‌ی دیده‌نشده معنای کلمات را می‌یابد. در [20] روشی مبتنی بر شباهت کسینوسی ارائه شده که از دو ویژگی کلمات متواتر متن و کلمات اطراف کلمه‌ی مبهم برای ابهام‌زدایی استفاده می‌کند.

• روش بدون ناظر<sup>۱۱</sup>

برخلاف روش‌های با ناظر، این روش به مجموعه‌داده‌ی برچسب‌دار نیازی ندارد و بر مبنای این کار می‌کند که کلمات با معانی یکسان، در یک محتوای یکسان ظاهر می‌شوند [11]. چونهی ژانگ [21]

<sup>1</sup> Simplified Lesk

<sup>2</sup> Supervised

<sup>3</sup> Classifier

<sup>4</sup> Naïve Bayes

<sup>5</sup> Part of speech (POS)

<sup>6</sup> Collocation

<sup>7</sup> Syntactic Relations

<sup>8</sup> Support-vector machine

<sup>9</sup> Neural Network

<sup>10</sup> Decision list

<sup>11</sup> Unsupervised



روشی مبتنی بر ژنتیک الگوریتم به نام ابهام‌زدایی معنایی ژنتیکی<sup>۱</sup> ارائه داده است که ابتدا از وردنت تمام هم‌نشیم‌های مورد نظر را استخراج و سپس از الگوریتم ژنتیک ساده و وزن‌دار برای بیشینه‌کردن میزان شباهت معنایی استفاده می‌کند. پدرس‌ن در [22] با استفاده از برداری از ویژگی‌ها برای ابهام‌زدایی، خوشه‌بندی<sup>۲</sup> ارائه داده است. ورونیس در [23] برای برنامه‌های بازیابی اطلاعات<sup>۳</sup> یک الگوریتم بدون ناظر مبتنی بر گراف برای ابهام‌زدایی معنایی به نام هایپرلکس<sup>۴</sup> ارائه داده است.

لازم به ذکر است که با استفاده از الگوریتم بوت‌استرپ<sup>۵</sup> و یادگیری ترارسانی<sup>۶</sup> می‌توان الگوریتم‌هایی شبه با ناظر برای ابهام‌زدایی معنایی پیاده‌سازی نمود که میزان کمی داده‌ی برچسب‌گذاری شده نیاز دارند [24]. از میان روش‌های ارائه‌شده در زبان فارسی می‌توان به دو مطالعه‌ی زیر اشاره کرد:

- روش ارائه شده در [25] گراف وابستگی معنایی<sup>۷</sup> است، که گره‌ها معانی مختلف کلمات درون یک جمله‌اند و با دنبال کردن یک مسیر معانی برای هر کلمه یک معنا انتخاب می‌شود. در این روش وزن هر گره توسط الگوریتم مرکزیت<sup>۸</sup> و وزن یال‌ها از سه روش جی-پی‌ام‌آی<sup>۹</sup>، پی‌ام‌آی<sup>۱۰</sup> و اس-پی‌ام‌آی<sup>۱۱</sup> محاسبه می‌شوند. معنای هر کلمه توسط دو روش جی‌ای‌دبلیو<sup>۱۲</sup> که مسیری با بیش‌ترین جمع وزن یال‌ها و جی‌ان‌دبلیو<sup>۱۳</sup> که گره‌های با بیش‌ترین وزن را به عنوان معنی انتخاب می‌کنند، محاسبه شده است.
- در [26] دو روش با ناظر کی‌ان‌ان<sup>۱۴</sup> و بیز ساده مقایسه شده‌اند. برای این دو روش مجموعه‌داده‌ی دارای برچسب نیاز بوده‌است، که به صورت دستی مجموعه‌داده‌ی همشهری برچسب‌گذاری شده‌است. معیار شباهت کسینوسی نرمال‌شده در نظر گرفته شده و در طبقه‌بندها از برداری از ویژگی‌ها استفاده شده است.

<sup>1</sup> genetic word sense disambiguation (GWSD)

<sup>2</sup> Clustering

<sup>3</sup> Information Retrieval applications

<sup>4</sup> HyperLex

<sup>5</sup> Bootstrap

<sup>6</sup> Transductive learning

<sup>7</sup> Semantic Dependency graph

<sup>8</sup> Centrality

<sup>9</sup> J-PMI

<sup>10</sup> PMI

<sup>11</sup> S-PMI

<sup>12</sup> GEW

<sup>13</sup> GNW

<sup>14</sup> kNN

## ۲-۳-۲- پیوندهای موجودیت

همان‌طور که در بخش ۱ بیان کردیم، پیوندهای موجودیت، نگاشت کردن موجودیت‌های موجود در یک متن به موجودیت‌های یک پایگاه‌دانش است. گلو<sup>۱</sup> یک سیستم پیوندهای موجودیت است، که توسط راتینوو در [27] ارائه شده است. این سیستم از پایگاه‌دانش ویکی‌پدیا استفاده می‌کند و ابتدا چند کاندید برای یک موجودیت انتخاب و آن‌ها را رتبه‌بندی می‌کند. سپس تصمیم می‌گیرد که آیا موجودیت با بالاترین رتبه می‌بایست به موجودیت مورد نظر نگاشت شود یا نگاشتی برای این موجودیت وجود ندارد. هوفارت در [28] سیستمی به نام آیدا<sup>۲</sup> ارائه داده است که با ساخت گرافی جدید از پایگاه‌دانش و با امتیازدهی با استفاده از احتمال و شباهت کاندیدها عمل پیوندهای را انجام می‌دهد. در مطالعه‌ای دیگر گوین به همراه هوفارت و همکاران [29] نسخه‌ای دیگر از آیدا به نام آیدا-لایت<sup>۳</sup> ارائه داده‌اند. هر دو آیدا و آیدا-لایت از یاگو<sup>۴</sup> به عنوان پایگاه‌دانش استفاده می‌کنند. آیدا-لایت از یک روش دولایه برای ابهام‌زدایی پیوندهای استفاده می‌کند. به این گونه که ابتدا موجودیت‌هایی که ابهام‌زدایی پیوندهای آن‌ها کم هزینه‌تر و ساده‌تر است، پیوندهای شده و سپس با داشتن این پیوندهای و محتوای مربوطه در لایه بعد به ابهام‌زدایی پیوندهای پر هزینه‌تر و سخت‌تر پرداخته می‌شود. پی‌پی‌آرسیم<sup>۵</sup> سیستم ارائه‌شده در [30] است، که از تلفیق اطلاعات محلی و سراسری استفاده می‌کند و الگوریتمی بر پایه‌ی گراف برای پیوندهای موجودیت است. در این سیستم از سرتیتر ویکی‌پدیا برای تولید کاندیدها استفاده می‌شود. پیوندهای و شناخت مشترک موجودیت‌های نام‌دار<sup>۶</sup> که در [31] ارائه شده است، از اولین مطالعاتی است که شناخت و پیوندهای موجودیت نام‌دار را به عنوان یک کار مشترک انجام داده است. در [32] از سی‌ان‌ان<sup>۷</sup> برای استخراج اطلاعات در راستای پیوندهای موجودیت استفاده شده است. در [33] از روش جاگذاری<sup>۸</sup> به منظور ابهام‌زدایی پیوندهای موجودیت استفاده شده است. در [10] روشی برای پیوندهای موجودیت فارسی با استفاده از فارس‌بیس ارائه داده شده است. در این روش ابتدا تشابه کسینوسی میان محتوای متن و محتوای متنی مقاله‌ی ویکی‌پدیا محاسبه و با استفاده از فرمولی با شباهت بین محتوای متن و ابرپیوندهای مقاله ترکیب و امتیاز نهایی هر کاندید را مشخص می‌کنند.

<sup>1</sup> GLOW

<sup>2</sup> AIDA

<sup>3</sup> AIDA-light

<sup>4</sup> YAGO

<sup>5</sup> PPRSim

<sup>6</sup> Joint named entity recognition and disambiguation (JERL)

<sup>7</sup> Convolutional neural networks (CNN)

<sup>8</sup> Embedding method

## ۲-۳-۳- نگاشات میان ابهام‌زدایی معنایی و پیونددهی موجودیت

ناویگلی و همکاران ابتدا در [34] پایگاه‌دانش ویکی‌پدیا را به وردنت انگلیسی نگاشت کردند و بیبل‌نت<sup>۱</sup> را ارائه‌دادند و سپس، در [35] الگوریتمی به نام ببیلیفای<sup>۲</sup> برای ابهام‌زدایی معرفی کردند. در نگاشت ارائه‌شده از هم‌معنی‌ها، هایپر‌نیم‌ها، هیپونیم‌ها، خواهر و گلاس مربوطه برای هر هم‌نشیم وردنت و برچسب داخل پرانتز، لینک‌ها و دسته‌بندی‌های ویکی‌پدیا استفاده می‌شود.

## ۲-۴- نتیجه‌گیری

در این بخش مطالعات پیشین مربوط به حوزه‌ی مدنظر بررسی شدند. همان‌گونه که مطالعه شد، با وجود پیشرفت ابهام‌زدایی معنایی و پیونددهی موجودیت به صورت جداگانه و وجود نگاشت میان پایگاه‌دانش و وردنت انگلیسی، برای زبان فارسی این نگاشت وجود نداشته و نیاز به آن در این حوزه احساس می‌شود. در مطالعه‌ی پیش رو سعی در نگاشت فارسن‌نت و فارس‌بیس شده است.

---

<sup>1</sup> BabelNet

<sup>2</sup> Babelify

## فصل ۳ :

### روش تحقیق

### ۳-۱- مقدمه

در این بخش به بررسی روش ارائه شده در این پروژه برای نگاشت فارسنت به فارس بیس می پردازیم.

### ۳-۲- تشریح کامل روش تحقیق

در ابتدا قصد بر انجام روشی همانند روش ارائه شده برای ساخت بیبلنت [34] کردیم. به این معنا که از فارس بیس موارد زیر:

- سرتیترهای داخل پرانتز در صورت موجود
- پیوندها
- دسته بندی ها

و از فارسنت موارد زیر را جمع آوری کردیم.

- کلمات درون یک هم نشیم
- هاپرنیم و پیپونیم ها
- خواهرها
- گلاس مربوطه

سپس با استفاده از فرمول زیر که در مقاله ی بیبلنت آورده شده بود امتیاز هر هم نشیم را محاسبه می کنیم.

(۳-۱)

$$score(s, w) = |Ctx(s) \cap Ctx(w)| + 1$$

$$p(s, w) = \frac{score(s, w)}{\sum_{s0 \in SensesWN(w)} \sum_{w0 \in SensesWiki(w)} score(s0, w0)} \quad (۳-۲)$$

اما به دلیل پایین بودن دقت این روش تصمیم به تعویض روش و یافتن روشی جدید برای نگاشت فارسنت به فارس گرفتیم.

در روش جدید بنا بر آن شد که ابتدا با استفاده از کد [10] الگوریتم پیونددهی موجودیت معرفی شده در

این مقاله را بر روی مجموعه داده اجرا کرده و سپس الگوریتم ابهام زدایی معنایی مشابه با این مقاله پیاده سازی کرده و روی مجموعه داده اجرا کنیم.

در این روش برای هر جمله ی ورودی ابتدا اجزای نحوی آن و سپس هم نشیم های کاندید برای هر کلمه پس از پیش پردازش استخراج می شوند. سپس هم نشیم هایی که با اجزای نحوی هم خوانی ندارند حذف شده و برای هم نشیم های باقی مانده موارد زیر استخراج می شوند.

- کلمات درون یک هم نشیم
- هاپرنیم و پیپونیم ها
- خواهرها
- گلاس مربوطه
- مثال مربوطه

پس از بدست آوردن این موارد امتیاز با استفاده از جمع تشابهات میان کلمات درون این موارد و کلمات محتوای مربوطه بدست می آید و ابهام زدایی معنایی انجام می شود.

پس از انجام ابهام زدایی معنایی بر روی مجموعه داده ی مورد نظر، نتایج پیونددهی و ابهام زدایی با یکدیگر مقایسه شده و در صورت نگاشت به بخش مشترک به یکدیگر نگاشت می شوند و تواتر این نگاشت ها محاسبه می شود. در آخر با استفاده از یک مجموعه داده ی برچسب دار از نگاشت بین فارسی و فارسی بیس یک آستانه ی مناسب برای قبول و یا عدم قبول نگاشت انتخاب می شود.

## فصل ۴ :

### نتایج و تفسیر آنها

#### ۴-۱- مقدمه

در این بخش نتایج حاصل از نگاشت بین فارس‌نت و فارس‌بیس با استفاده از روش مطرح شده در بخش ۳ را بررسی می‌کنیم.

#### ۴-۲- محتوا

##### ۴-۲-۱- مجموعه داده‌ها

به منظور ارزیابی الگوریتم ارائه شده و به دلیل جدید بودن نگاشت بین فارس‌نت و فارس‌بیس مجموعه داده‌ای برای ارزیابی و یا الگوریتمی برای مقایسه‌ی نتایج وجود نداشت. به همین دلیل، از مجموعه داده‌ای که توسط یک شخص خبره تولید و شامل ۷۶۷ نگاشت بین فارس‌نت و فارس‌بیس بود، برای ارزیابی استفاده کردیم. همچنین، الگوریتم‌های ابهام‌زدایی معنایی و پیونددهی موجودیت بر روی مجموعه داده‌ی همشهری اجرا شده‌اند.

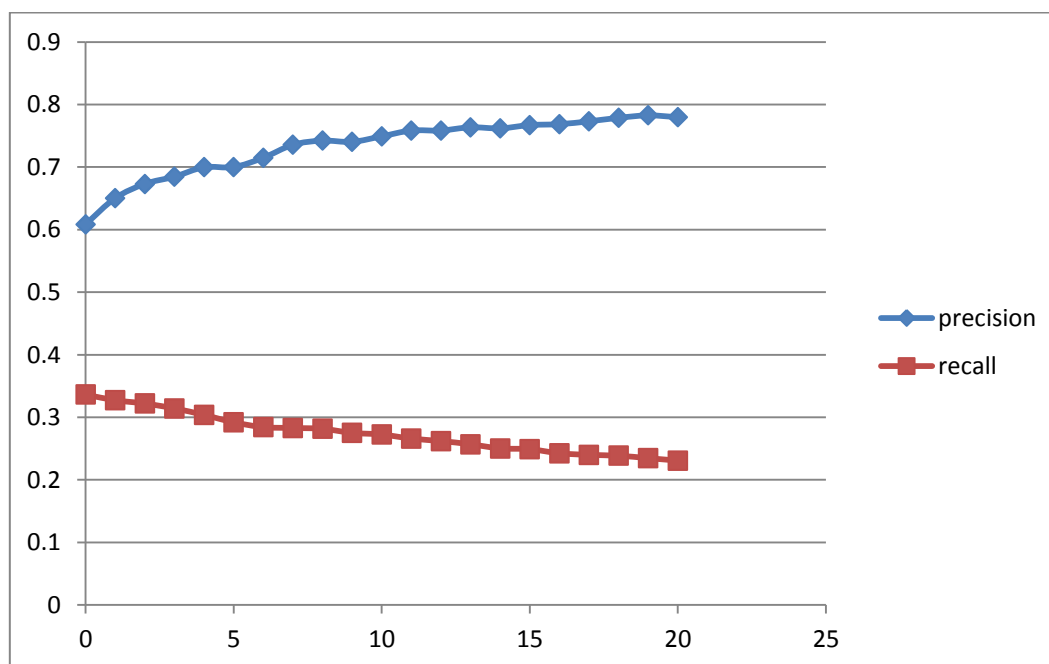
##### ۴-۲-۲- ارزیابی نتایج

نتایج به دست آمده از روش ارائه شده را پس از اعمال آستانه‌های مختلف برای قبولی با مجموعه داده‌ی ارزیابی مقایسه کردیم. این مقایسه با دو معیار دقت و بازخوانی و به منظور یافتن آستانه‌ی بهینه انجام شده است.

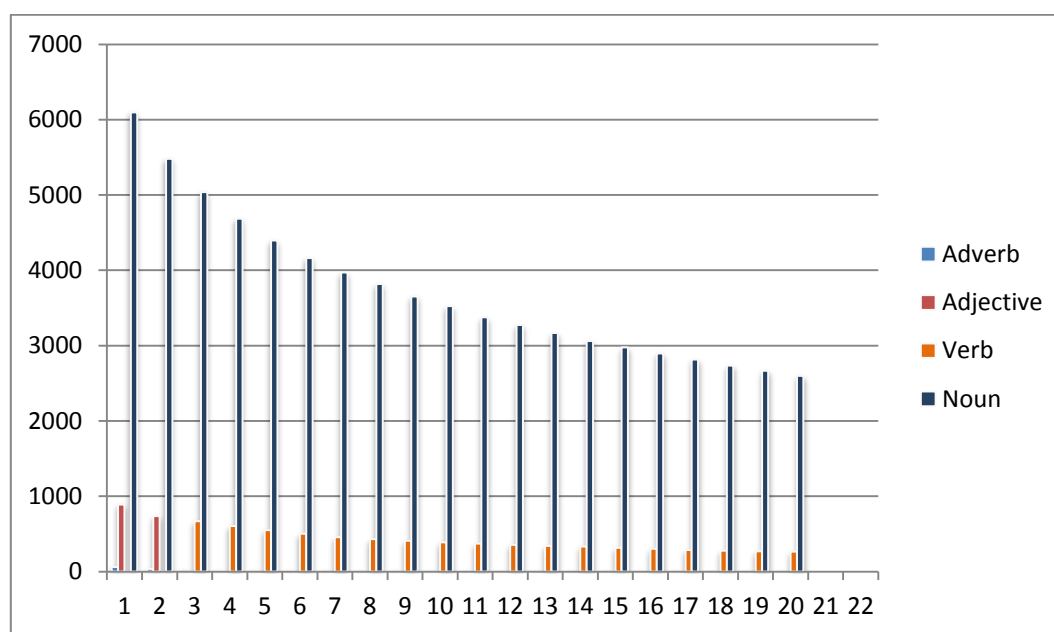
شکل (۱-۴) افزایش و کاهش معیار دقت و بازخوانی با افزایش آستانه را نمایش می‌دهد. همان گونه که قابل انتظار است، با افزایش آستانه‌ی تواتر نگاشت قابل قبول معیار دقت افزایش و معیار بازخوانی کاهش پیدا می‌کند. همچنین، همان طور که در شکل (۳-۴) قابل مشاهده است، افزایش آستانه تعداد نگاشت‌ها را کاهش می‌دهد. در شکل (۲-۴) تعداد هر یک از اجزای سخن تشخیص داده شده در این نگاشت را مشاهده می‌کنیم. تعداد اسم‌های تشخیص داده شده از سایر اجزای سخن بیش‌تر است و بعد از اسم‌ها، فعل‌ها



بیشترین‌اند.

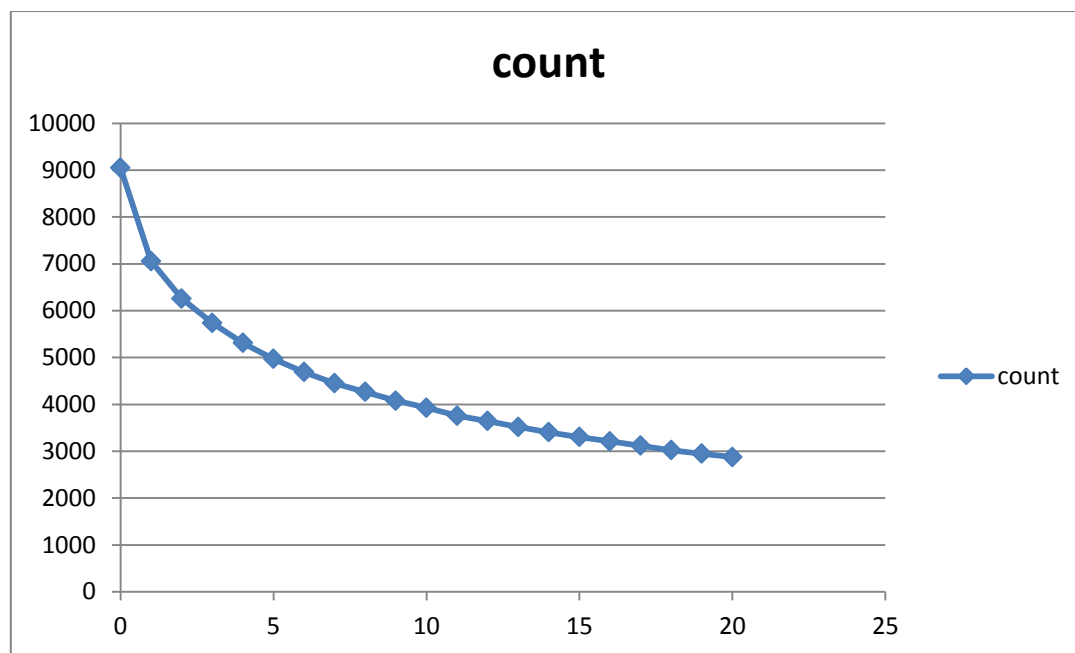


شکل (۴-۱) مقدار معیار دقت و بازخوانی برای آستانه‌های متفاوت



شکل (۴-۲) تعداد اجزای سخن در میان نگاشت‌ها برای هر آستانه

در شکل (۴-۳) تعداد نگاشت‌ها برای هر آستانه را نشان می‌دهد. برای تواتر حداقل یکبار تعداد ۹۰۰۰ نگاشت و برای آستانه‌ی حداقل ۲۰ تعداد ۳۰۰۰ نگاشت انجام شده است.



شکل (۳-۴) تعداد نگاشت‌ها در هر آستانه

با توجه به سه نمودار نمایش داده شده، آستانه‌ی ۵ به نظر مناسب‌تر از بقیه‌ی آستانه‌ها می‌آید. با گرفتن آستانه‌ی ۵، معیار دقت برابر ۷۰، بازخوانی برابر ۳۰ و تعداد نگاشت‌ها برابر ۵۰۰۰ است.

## فصل ۵ :

### جمع‌بندی و پیشنهادها

## ۵-۱- مقدمه

ابهام‌زدایی معنایی و پیونددهی موجودیت از موضوعات مطرح در حوزه‌ی پردازش زبان طبیعی‌اند، که با استفاده از فارس‌نت و فارس‌بیس انجام می‌گیرند. در این مطالعه سعی در پیوند این دو حوزه و نگاشت بین فارس‌نت و فارس‌بیس شده‌است. در این بخش به معرفی کلی مطالعه‌ی انجام شده و جمع‌بندی می‌پردازیم. در انتها نیز درباره‌ی نوآوری و پیشنهادات آتی برای این مطالعه بحث می‌کنیم.

## ۵-۲- جمع‌بندی

در این مطالعه، سعی در نگاشت فارس‌نت و فارس‌بیس که تاکنون برای زبان فارسی انجام نشده‌است، کرده‌ایم. برای این کار، ابتدا یک الگوریتم پیونددهی موجودیت مناسب یافته و سپس یک الگوریتم برای ابهام‌زدایی معنایی مشابه آن پیاده‌سازی نمودیم. سپس اشتراک خروجی این دو الگوریتم را گرفته و تواتر آن‌ها را حساب کردیم. در آخر با استفاده از معیار دقت، بازخوانی و تعداد نگاشت‌ها آستانه‌ای برای نگاشت با تواتر قابل قبول در نظر گرفتیم. معیار دقت، بازخوانی و تعداد نگاشت نهایی برابر با ۷۰ درصد، ۳۰ درصد و ۵۰۰۰ نگاشت بوده‌است.

## ۵-۳- نوآوری

نگاشت میان فارس‌نت و فارس‌بیس به منظور اجرای هم‌زمان و بهینه‌ی ابهام‌زدایی معنایی و پیونددهی موجودیت برای زبان فارسی تاکنون انجام نشده‌است. در این مطالعه، سعی در نگاشت این دو منبع مهم در حوزه‌ی پردازش زبان طبیعی کرده‌ایم، که خود نوآوری محسوب می‌شود.

## ۵-۴- پیشنهادها

برای کارهای آتی می‌توان الگوریتم را بر روی مجموعه‌داده‌های متنوع و بیش‌تری اجرا کرد تا نگاشت‌های

فعل‌ها و تواتر نگاشت‌ها بهبود پیدا کند. همچنین می‌توان با نگاشت‌های یافت شده، یک الگوریتم برای ابهام‌زدایی و پیونددهی موجودیت هم‌زمان، همانند بیبل‌نت [34]، ارائه داد.

## مراجع

- [1] J. Sreedhar, S. Raju, A. Babu, A. Shaik, and P. Kumar, "Word Sense Disambiguation: An Empirical Survey," *International Journal of Soft Computing and Engineering*, Volume-2, Issue-2, pp. 494-503, May 2012
- [2] A. Moro, A. Raganato, and R. Navigli, "Entity Linking meets Word Sense Disambiguation: a Unified Approach," *Trans. Assoc. Comput. Linguist.*, vol. 2, pp. 231–244, 2018.
- [3] A. Moro, F. Cecconi, and R. Navigli, "Multilingual word sense disambiguation and entity linking for everybody," *CEUR Workshop Proc.*, vol. 1272, pp. 25–28, 2014.
- [4] S. Kumar, N. Sharma, and S. Niranjana, "Word Sense Disambiguation Using Association Rules : A Survey," *Int. J. Comput. Technol. Electron. Eng.*, vol. 2, no. 2, pp. 93–98, 2012.
- [5] B. Hachey, W. Radford, and J. R. Curran, "Graph-based named entity linking with Wikipedia," *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 6997 LNCS, pp. 213–226, 2011.
- [6] T. Huber and P. Definition, "Entity Linking - A Survey of Recent Approaches," *Department of Computer Science, Humboldt-Universität zu Berlin* pp. 1–5, 2012.
- [7] M. Shamsfard et al., "Semi automatic development of farsnet; the persian wordnet," *Proc. 5th Glob. WordNet Conf. Mumbai, India*, vol. 29, 2010.
- [8] A. M. George, "WordNet: A Lexical Database for English," *Commun. ACM*, vol. 38, pp. 39–41, 1995.
- [9] M. B. Sajadi, B. M. Bidgoli, and A. Hadian, "FarsBase: A cross-domain farsi knowledge graph," *CEUR Workshop Proc.*, vol. 2198, 2018.
- [10] M. Asgari, A. Hadian, and B. Minaei-Bidgoli, "FarsBase: The Persian Knowledge Graph," *Semantic Web J.*, vol. 1, no. 0, pp. 1–5, 2018.
- [11] P. P. Borah, G. Talukdar, and A. Baruah, "Approaches for Word Sense Disambiguation – A Survey," *Int. J. Recent Technol. Eng.*, vol. 3, no. 1, pp. 35–38, 2014.
- [12] M. Lesk, "Automatic Sense Disambiguation Using Machine Readable Dictionaries: How to Tell a Pine Cone from an Ice Cream Cone," *Proc. 5th Annu. Int. Conf. Syst. Doc.*, pp. 24–26, 1986.
- [13] A. Kilgarriff, "{E}nglish {SENSEVAL}: Report and Results," *Proc. Int. Conf. Lang. Resour. Eval.*, 1998.
- [14] W. A. Gale, K. W. Church, and D. Yarowsky, "A method for disambiguating word senses in a large corpus," *Comput. Hum.*, vol. 26, no. 5–6, pp. 415–439, 1992.
- [15] Y. K. Lee, H. T. Ng, and T. K. Chia, "Supervised word sense disambiguation with support vector

- machines and multiple knowledge sources.” Senseval-3 Third Int. Work. Eval. Syst. Semant. Anal. Text, no. July, pp. 137–140, 2004.
- [16] R. Bruce and J. Wiebe, “Word-sense disambiguation using decomposable models,” In 32nd Annual Meeting of the Association for Computational Linguistics, Las Cruces, NM, pp. 139–146, 2007.
- [17] G. Towell, “Disambiguating Highly Ambiguous Words,” In: Computational Linguistics, March 1998, Vol. 24, Number 1, vol. 20899, pp. 1–22, 2002.
- [18] D. Yarowsky, “Decision Lists for Lexical Ambiguity Resolution: Application to Accent Restoration in Spanish and French,” Natural Language Engineering, 2002.
- [19] N. Chatterjee and R. Misra, “Word-Sense Disambiguation using maximum entropy model,” in 2009 Proceeding of International Conference on Methods and Models in Computer Science (ICM2CS).
- [20] M. Nameh, S. M. Fakhrahmad, and M. Z. Jahromi, “A New Approach to Word Sense Disambiguation Based on Context Similarity,” Proc. World Congr. Eng., vol. I, pp. 8–11, 2011.
- [21] C. H. Zhang, Y. Zhou, and T. Martin, “Genetic word sense disambiguation algorithm,” Proc. - 2008 2nd Int. Symp. Intell. Inf. Technol. Appl. IITA 2008, vol. 1, pp. 123–127, 2008.
- [22] T. Pedersen and R. Bruce, “Unsupervised Text Mining,” Dallas TX, Departement of computer science and engineering, southern methodist university 1997.
- [23] J. Véronis, “HyperLex: Lexical cartography for information retrieval,” Comput. Speech Lang., vol. 18, no. 3 SPEC. ISS., pp. 223–252, 2004.
- [24] T. P. Pham, H. T. Ng, and W. S. Lee, “Word sense disambiguation with semisupervised learning,” Proc. 20th Natl. Conf. Artif. Intell. - Vol. 3, pp. 1093–1098, 2005.
- [25] M. Soltani and H. Faili, “A statistical approach on Persian word sense disambiguation,” 7th Int. Conf. Informatics Syst., pp. 1–6, 2010.
- [26] M. Hamidi, A. Borji, and S. S. Ghidary, “Persian Word Sense Disambiguation,” 15th Iran. Conf. Electr. Electron. Eng. Proc., pp. 114–118, 2007.
- [27] L. Ratnov, D. Roth, D. Downey, and M. Anderson, “Local and Global Algorithms for Disambiguation to Wikipedia,” Annu. Meet. Assoc. Comput. Linguist., pp. 1375–1384, 2011.
- [28] J. Hoffart et al., “Robust Disambiguation of Named Entities in Text Johannes,” Proc. 2011 Conf. Empir. Methods Nat. Lang., pp. 782–792, 2011.
- [29] D. B. Nguyen, J. Hoffart, M. Theobald, and G. Weikum, “AIDA-light: High-throughput named-entity disambiguation,” CEUR Workshop Proc., vol. 1184, 2014.
- [30] M. Pershina, Y. He, and R. Grishman, “Personalized Page Rank for Named Entity Disambiguation,” Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 238–243, 2015.



- [31] T. Ishizaki, K. Yokochi, K. Chiba, T. Tabuchi, and T. Wagatsuma, "Joint Named Entity Recognition and Disambiguation," *Pediatr. Pharmacol. (New York)*, vol. 1, no. 4, pp. 291–303, 1981.
- [32] M. Francis-Landau, G. Durrett, and D. Klein, "Capturing Semantic Similarity for Entity Linking with Convolutional Neural Networks," *arXiv preprint arXiv:1604.00734*, 2016.
- [33] I. Yamada, H. Shindo, H. Takeda, and Y. Takefuji, "Joint Learning of the Embedding of Words and Entities for Named Entity Disambiguation," *The SIGNLL Conference on Computational Natural Language Learning (CoNLL)* 2016.
- [34] R. Navigli and S. Ponzetto, "BabelNet: Building a very large multilingual semantic network," *Proc. 48th Annu. Meet. ...*, no. July, pp. 216–225, 2010.
- [35] T. Flati and R. Navigli, "Three birds (in the LLOD cloud) with one stone: BabelNet, Babelfy and the Wikipedia Bitaxonomy," *Proc. Semant.*, pp. 4–5, 2014.

## فهرست واژگان فارسی به انگلیسی

Threshold	آستانه
Word Sense Disambiguation (WSD)	ابهام‌زدایی معنایی
Genetic word sense disambiguation (GWSD)	ابهام‌زدایی معنایی ژنتیکی
Part of speech (POS)	اجزای کلام
Supervised	با ناظر
Unsupervised	بی‌ناظر
Information Retrieval applications	برنامه‌های بازیابی اطلاعات
Naïve Bayes	بیز ساده
Natural Language Processing (NLP)	پردازش زبان طبیعی
Named-Entity Linking	پیونددهی موجودیت
Antonymy	تضاد
Clustering	خوشه‌بندی
Gold dataset	دادگان طلایی
Syntactic Relations	روابط نحوی
Embedding method	روش جایگذاری
Neural Network	شبکه‌ی عصبی
Classifier	رده‌بند
Semantic Dependency graph	گراف وابستگی معنایی
Decision list	لیست تصمیم
Support-vector machine	ماشین بردار پشتیبانی
Knowledge based	مبتنی بر دانش
Centrality	مرکزیت
Collocation	همایند
Synsets	هم‌نشیم‌ها
Transductive learning	یادگیری ترانسدوکتی

## فهرست واژگان انگلیسی به فارسی

Antonymy	تضاد
Centrality	مرکزیت
Classifier	رده‌بند
Clustering	خوشه‌بندی
Collocation	همایند
Decision list	لیست تصمیم
Embedding method	روش جایگذاری
Genetic word sense disambiguation (GWSD)	ابهام‌زدایی معنایی ژنتیک
Gold dataset	دادگان طلایی
Information Retrieval applications	برنامه‌های بازیابی اطلاعات
Knowledge based	مبتنی بر دانش
Named-Entity Linking	پیونددهی موجودیت
Natural Language Processing (NLP)	پردازش زبان طبیعی
Naïve Bayes	بیز ساده
Neural Network	شبکه‌ی عصبی
Part of speech (POS)	اجزای کلام
Semantic Dependency graph	گراف وابستگی معنایی
Supervised	با ناظر
Support-vector machine	ماشین بردار پشتیبانی
Synsets	هم‌نشیم‌ها
Syntactic Relations	روابط نحوی
Transductive learning	یادگیری ترانسانی
Threshold	آستانه
Unsupervised	بی‌ناظر
Word Sense Disambiguation (WSD)	ابهام‌زدایی معنایی

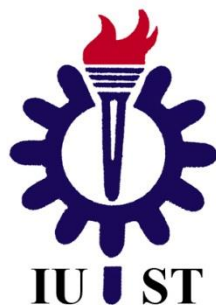
## **Abstract:**

Word sense disambiguation, which maps each word to a Synset of a Wordnet, is one important problem in the field of text knowledge extraction. This problem is related to the problem of Entity linking, which links words to corresponding entities of a knowledge graph.

In Persian, these two problems can be solved using Persian Wordnet (Farsnet) and Persian knowledge graph (FarsBase). Although there exists a uniform combination of these two problems for English language, there are no such work for Persian language.

The main aim of this project is to represent a mapping between Farsnet and Farsbase. First of all, our proposed method applies a Persian entity linker on a corpus and then, uses a word sense disambiguation algorithm similar to the entity linker for the same corpus. Then, we check the result of the entity linker and the word sense disambiguation algorithm for similarities, in order to map linked entities to word senses (Synsets of FarsNet). At the end, we count the number of occurrence of each mapped pairs and look for an appropriate threshold for accepting a mapped pair. The final mapping is evaluated using a dataset of 767 manually mapped entity-Synset pairs. The result had a precision and recall of 70 and 30 presents and contains 5000 mapped pairs.

**Keywords: Word Sense Disambiguation- Entity Linking- FarsNet- FarsBase**



**Iran University of Science and Technology**  
**Computer engineering Department**

# **Entity linking algorithm for Word sense disambiguation of Persian language**

**A Thesis Submitted in Partial Fulfillment of the Requirement for the  
Degree of Bachelors in Computer engineering**

**By:**  
**Banafsheh Karimian**

**Supervisor:**  
**Dr. Minaei**

**June 2019**