

## بررسی انواع منابع دارای ساختار و منابع بدون ساختار و پیش پردازش های ابهام زدایی مفهوم کلمات در پردازش زبان طبیعی

علی وحید رودسری<sup>۱</sup>، کبری خوشرفتار<sup>۲</sup>، مرتضی گل پور<sup>۳</sup>، اعظم عندلیب<sup>۴</sup>

۱- کارشناسی ارشد مهندسی فناوری اطلاعات، دانشگاه گیلان، گروه کامپیوتر، رشت، ایران

۲- دانشجوی کارشناسی ارشد مهندسی نرم افزار دانشگاه آزاد اسلامی واحد چالوس،

گروه کامپیوتر، مازندران، ایران

۳- دانشجوی کارشناسی ارشد مهندسی نرم افزار، دانشگاه آزاد اسلامی واحد چالوس،

گروه کامپیوتر، مازندران

۴- عضو هیات علمی گروه مهندسی کامپیوتر، واحد رشت، دانشگاه آزاد اسلامی، رشت، گیلان، ایران

### چکیده

دانش مهمترین بخش ابهام زدایی مفهوم کلمات است. این دانش ها می توانند در شکل های گوناگون و به صورت یک مجموعه از متون باشند که در آن مفهوم کلمه برچسب گذاری شده است. پایگاه دانش یک مجموعه از متن، برچسب ها و توضیحات در جهت تشخیص مفهوم کلمه است. مانند فرهنگ لغت قابل خواندن توسط ماشین، شبکه معنایی، اصطلاحنامه و آنتولوژی. تقریباً از تمام این منابع در ابهام زدایی مفهوم کلمات استفاده می شود. کلیه منابع به دو دسته منابع دارای ساختار و منابع بدون ساختار تقسیم می شوند. جمله ورودی، یک متن بدون ساختار از اطلاعات است. برای کسب مفهوم صحیح کلمات باید بر روی آن پیش پردازش هایی انجام شود تا بتوانیم بستری را فراهم نماییم که بتوان بهترین مفهوم را بدست آورد. در این مقاله، منابع دارای ساختار و منابع بدون ساختار و پیش پردازش های ابهام زدایی مفهوم کلمات در پردازش زبان طبیعی را مورد بررسی قرار می دهیم که بر اساس بررسی های انجام شده، استفاده از وردنت پیشنهاد می شود که یک منبع ضروری برای ابهام زدایی مفهوم کلمات است و یک منبع دارای ساختار می باشد.

**کلمات کلیدی:** اصطلاحنامه، فرهنگ لغت های قابل خواندن توسط ماشین، آنتولوژی، وردنت، Corpora.

## ۱- مقدمه

یک مسئله ی اصلی در پردازش زبان طبیعی، ابهام زدایی مفهوم کلمات است. وظیفه ی گسترده ای از روش ها برای ابهام زدایی مفهوم کلمات استفاده شده است. با این حال، اغلب به دانش های دستی مانند: فرهنگ لغت قابل خواندن توسط ماشین یا واژه فهرست، شبکه های معنایی، نیازمند است (Pedersen and Bruce, ۱۹۹۸). ابهام زدایی مفهوم کلمات وظیفه ی تعیین مفاهیم کلماتی را که دارای معانی زیادی در یک متن خاص هستند را برعهده دارد (Ho et al, ۲۰۱۰). با وجود اهمیت کار ابهام زدایی مفهوم کلمات و محبوبیت آن به عنوان یک موضوع مورد مطالعه، ابزار و منابعی که ابهام زدایی مفهوم کلمات حمایت کنند، تعمیم و استاندارد نسبتاً کمی دیده می شود (Miller et al, ۲۰۱۳). ابهام زدایی مفهوم کلمات به طور معمول به عنوان یک مسئله ی طبقه بندی نشان داده می شود که در آن، به هر کلمه ی مبهم یک برجسب مفهومی از فهرست مفاهیم از پیش تعریف شده که در طول فرآیند ابهام زدایی مفهوم کلمات وجود دارد اختصاص داده شده است. در اولین گام برای ابهام زدایی مفهوم کلمات باید عبارت را به یک عبارت رسمی براساس قواعد نحوی تعریف شده تبدیل کرد و سپس براساس یک پایگاه دانش معنی درست آن کلمه را بدست آورد. در واقع اسکلت اصلی هر سیستم ابهام زدایی دارای دو قسمت است: پایگاه دانش و الگوریتمی که بتواند معنی صحیح را از پایگاه دانش استخراج نماید. یکی از موانع اصلی برای کارایی بالا در ابهام زدایی مفهوم کلمات، تنگنا و محدودیت در کسب دانش است. در سال های اخیر، دو روش اصلی در ابهام زدایی مفهوم کلمات مورد مطالعه قرار گرفته است. مانند روش با ناظر و روش های مبتنی بر دانش. بدون داشتن یک پایگاه دانش پیدا کردن معنی کلمه برای انسان و کامپیوتر امری محال می باشد. ایجاد پایگاه دانش دشوار است زیرا در هر زمان تغییرات بسیاری در زبان در حال انجام است. این مشکلات، مسائل اصلی در ابهام زدایی مفهوم کلمات هستند (Kilgariff, ۲۰۰۲). منابع به دو دسته منابع دارای ساختار و منابع بدون ساختار تقسیم بندی می شوند که مورد بررسی قرار می گیرند.

## ۲- منابع دارای ساختار

در اینگونه منابع، در صورتی که مفاهیم کلمات شامل رابطه ای با یکدیگر باشند آنها نیز مشخص می شوند. این ساختارها به شرح زیر است:

- اصطلاحنامه
- فرهنگ لغت های قابل خواندن توسط ماشین
- آنتولوژی
- وردنت

## ۲-۱ اصطلاحنامه

این منبع اطلاعاتی درباره ی ارتباط بین کلمات فراهم می کند مانند مشابهت، متضاد و همچنین رابطه های دیگری که بین کلمات است را نیز نشان می دهد (Navigli, ۲۰۰۹).

## ۲-۲ فرهنگ لغت های قابل خواندن توسط ماشین

یکی از مشهور ترین منابعی است که برای پردازش زبان های طبیعی استفاده می شود. زمانی که این فرهنگ لغت مطرح شد به عنوان اولین فرهنگ لغت الکترونیکی معرفی گردید. این فرهنگ لغت الکترونیکی ترکیبی از فرهنگ لغت های CED، LDOCE، ODOE، OALD است. فرهنگ لغت های قابل خواندن توسط ماشین برای اولین بار در سال ۱۹۸۰ ساخته شد (Tan, ۲۰۱۳)

و از آن زمان برای تبدیل شدن یک منبع دانش برای مدل سازی زبان انسانی بکار می رود. اصطلاحنامه روابط وابسته به فرهنگ نویسی پایه را ارائه می دهد. نظیر مترادف و متضاد و احتمالا سایر روابط معنایی مانند: hypernymy، hyponymy، metronymy را فراهم می کند. فلپاوم<sup>۱</sup> در سال ۱۹۹۸ و میلر<sup>۲</sup> در سال ۱۹۹۰ اصطلاح نامه ی بین المللی روگت را که در سال ۱۸۵۲ توسط روگت<sup>۳</sup> ایجاد شده بود، کلمات را بر اساس روابط زیر طبقه بندی می کنند:

- مفاهیم انتزاعی
- فضا
- محتوا
- ایده های فکری
- اراده و تصمیم
- شهود اجتماعی و عاطفی

آخرین نسخه ی آن شامل ۲۵۰۰۰۰ موجودیت است. اصطلاحنامه ی مکواری که در سال ۱۹۸۷ توسط برنارد<sup>۴</sup> مورد بررسی قرار گرفت دارای بیش از ۲۰۰۰۰۰ مترادف مبتنی بر انگلیسی استرالیایی از جمله سبک گفتاری عامیانه ی استرالیایی است. شکل ۱ چهارچوب کلی را برای روش ابهام زدایی مفهوم کلمات مفهومی انطباقی که در این قسمت مطرح می شود را نشان می دهد. در اینجا فرآیند یادگیری توصیف می کند که با یک گام شامل کسب دانش از فرهنگ لغت های قابل خواندن توسط ماشین آغاز می شود. با استفاده از این کسب دانش، متن ورودی خوانده شده و یک گام ابهام زدایی آزمایشی انجام می گیرد. در مرحله ی تطابق، بیان می شود که ترکیبی از دانش اولیه بر اساس دانش جمع آوری شده از متن تا حدی ابهام زدایی می شود. دانش، با توجه به متن در دست تنظیم شده است و سپس برای نتیجه ی نهایی ابهام زدایی، دوباره به متن اعمال می شود. به عنوان

<sup>۱</sup> Fellbaum

<sup>۲</sup> Miller

<sup>۳</sup> Roget

<sup>۴</sup> Bernard

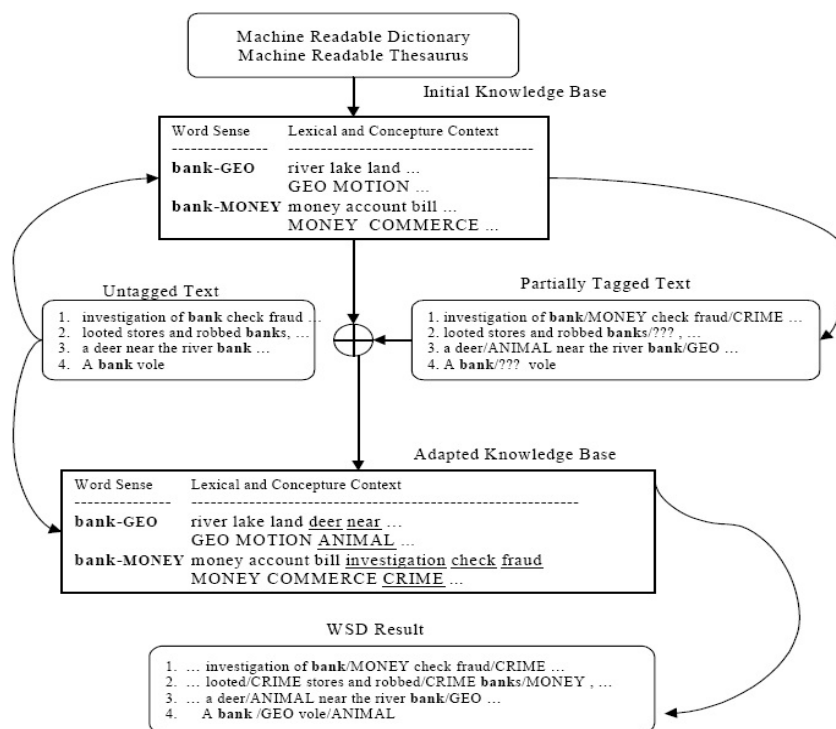
مثال: نمایش متنی اولیه از فرهنگ لغت انگلیسی معاصر لانگمن استخراج شده است. مفهوم bank\_GEO شامل هر دو اطلاعات لغوی و مفهومی است:

$\{\text{land, river, lake, } \dots\} \cup \{\text{GEO, MOTION, } \dots\}$

نمایش متنی اولیه، حاوی اطلاعات کافی برای ابهام زدایی عبارت زیر در متن ورودی است.

"A deer near the river bank"

در این گام ابهام زدایی آزمایشی، برچسب های مفهومی از deer/ANIMAL (آهو) و bank/GEO را تولید می کند. اما در بعضی از موارد bank به دلیل عدم آگاهی از دانش ابهام زدایی مفهوم کلمات، بدون برچسب هستند. ما مشاهده می کنیم که مفهوم bank\_GEO در زمینه Vole (موش صحرایی) از آنجایی که بین ANIMAL و GEOGRAPHY ارتباطی وجود ندارد، حل نشده است. پس از آن مرحله ای انطباق، deer (آهو) و ANIMAL را به نمایش متنی برای bank\_GEO اضافه می کند. تطابق نمایش متنی اولیه در حال حاضر با اطلاعات موجود بسیار غنی است و قادر به ابهام زدایی مثال bank در زمینه Vole (موش صحرایی) برای تولید نتیجه ی ابهام زدایی نهایی است (Chen, ۲۰۰۰).



## شکل ۱- ابهام زدایی مفهوم کلمات تطبیقی با استفاده از فرهنگ لغت های قابل خواندن توسط ماشین (Chen, ۲۰۰۰)

CR ها برای ۹ مفهوم اسمی از bank در فرهنگ لغت انگلیسی معاصر لانگمن (LDOCE) در زیر فهرست شده است. که در جدول ۱ نشان داده شده است:

- bank.۱.n.۱ land along the side of a river, lake, etc.;
- bank.۱.n.۲ earth which is heaped up in a field or garden, often making a border or division;
- bank.۱.n.۳ a mass of snow, clouds, mud, etc.;
- bank.۱.n.۴ a slope made at bends in a road or race-track, so that they are safer for cars to go round;
- bank.۱.n.۵ a high underwater of bank in a river, harbour, etc.;
- bank.۳.n.۱ a row, esp. of OAR in an ancient boat or KEY on a TYPEWRITER;
- bank.۴.n.۱ a place in which money is kept and paid out on demand, and where related activities go on;
- bank.۴.n.۲ a place where something is held ready for use, esp. ORGANIC products of human origin for medical use;
- bank.۴.n.۳ a person who keeps a supply of money or pieces for payment or use in a game of chance.

### جدول ۱- نمایش متنی واژگان برای مفهوم bank (Chen, ۲۰۰۰)

| Sense ID   | Sense Label <i>S</i> | Lexical Context Representation $LCR(D_{bank, s})$            |
|------------|----------------------|--|
| bank.4.n.1 | MONEY                | {place, money, keep, pay, demand, activity}                  |
| bank.1.n.1 | RIVER                | {land, lake, river}  |
| bank.1.n.5 | SANDBANK             | {underwater, sand, harbour}                                  |
| bank.1.n.2 | EARTH                | {earth, heap, field, garden, boarder, division}              |
| bank.1.n.3 | PILE                 | {mass, snow, cloud, mud}                                     |
| bank.1.n.4 | ROAD                 | {car, aircraft, move, side, turn}                            |
| bank.3.n.1 | ROW                  | {row, oar, boat, key, typewriter}                            |
| bank.4.n.2 | MEDICINE             | {place, hold, use, organic, product, human, origin, medical} |
| bank.4.n.3 | GAMBLE               | {person, keep, supply, money, payment, game, chance}         |

## ۳-۲ آنتولوژی

آنتولوژی (هستان شناسی)، یک مدل خاص از دامنه های ویژه است که عمدتاً شامل طبقه بندی و مجموعه ای از ارتباطات دارای

مفهوم است. داده هایی که در اینجا دارای رابطه هستند به صورت لینک هایی به هم وصل می شوند و هر یک از این لینک ها، یک نوع از رابطه را توصیف می کنند (Philpot et al, ۲۰۰۵). آنتولوژی، شامل بیش از ۸۰ هزار روابط بین مفهوم کلمات است. به عنوان مثال: abbreviation, hypernym, synonym (Tanaka et al, ۲۰۰۷). آنتولوژی مدلی انتزاعی از جهان واقع است که مفاهیم و روابط میان آن را در قلمروی مورد بحث نمایش می دهد. آنتولوژی ها که پایگاه دانش مفهومی هستند و در محدوده وسیعی از قلمروها کاربرد دارند که برای نمونه می توان به شبکه های جهان گستر معنایی، موتورهای جستجو، تجارت الکترونیک، پردازش زبان طبیعی، مهندسی دانش، استخراج و بازیابی اطلاعات، سیستم های چندعاملی، مدل سازی کیفی از سیستم های فیزیکی، طراحی پایگاه داده، سیستم های اطلاعات جغرافیایی و کتابخانه های رقمی اشاره نمود. در قلمروی کامپیوتر، آنتولوژی را می توان با یک چهارتایی (C,R,F,A) تعریف کرد که در آن:

- C مجموعه مفاهیم موجود در جهان مدل شده است.
- R مجموعه روابط میان مفاهیم است و خود به دو زیرمجموعه مجزای  $R_N$  و  $R_t$  افراز می شود:
- $R_t$  مجموعه روابط طبقه ای میان مفاهیم است که سلسله مراتب مشمول را ایجاد می کند و دودویی است.
- $R_N$  مجموعه روابط غیر طبقه ای است که ممکن است n تایی نیز باشد. ( $1 \leq n$ )
- F مجموعه تصریحات آنتولوژی در مورد مفاهیم و روابط آنها است و خود به دو زیرمجموعه  $F_t$  و  $F_N$  افراز می شود:
- $F_t$  مجموعه تصریحات آنتولوژی در باره روابط طبقه ای مفاهیم است. به عبارت دیگر، سلسله مراتب مشمول را نشان می دهد.
- $F_N$  مجموعه اصول بدیهی آنتولوژی درباره ی روابط غیر طبقه ای مفاهیم است.
- A مجموعه اصول بدیهی آنتولوژی است. که به زبان صوری، مثل منطق بیان می شود (شمس فرد و عبدالله زاده بارفروش، ۱۳۸۱).

تفاوت های اصلی بین آنتولوژی انسان و ماشین عبارتند از:

آنتولوژی انسان بسیار وسیع تر از آنتولوژی ماشین است. اما آنتولوژی ماشین باید رسمی باشد. یعنی با یک زبان قابل فهم توسط ماشین بیان شود. همچنین باید شفاف و دقیق باشد. یعنی تمام جزئیات به صورت کامل و غیر مبهم توصیف شده باشد. اما همان طور که ذکر شد، آنتولوژی انسان به صورت ناگهانی و تلویحی شکل می گیرد (نوروزی و طاهریان، ۱۳۹۰).

## ۲-۴ وردنت

وردنت، فرهنگ لغت الکترونیکی با قابلیت دسترسی آزاد است که شامل اسم، صفت و قید است و در دانشگاه پرینستون، توسعه داده شده است (Patwardhan et al, ۲۰۰۳). وردنت یک واژه نامه محاسباتی زبان انگلیسی در پردازش زبان طبیعی

است که به صورت گسترده مورد استفاده قرار می گیرد. وردنت یک منبع ضروری برای ابهام زدایی مفهوم کلمات است و وردنت یک منبع دارای ساختار است. در وردنت، همه ی مفاهیم مترادف کلمات در یک مجموعه قرار دارد که آن مجموعه ی مترادف نامیده می شود. آخرین نسخه ی این فرهنگ لغت ۳,۰ است و دارای ۱۵۵,۰۰۰ کلمه است که این کلمات در ۱۱۷,۰۰۰ گروه معنایی قرار گرفته اند. برای هر مجموعه مترادف یک مجموعه اطلاعات وجود دارد یعنی رابطه هایی وجود دارد که اطلاعاتی را برای کامل کردن این داده در اختیار کاربرها قرار می دهد.

به عنوان مثال در اینجا وقتی کلمه automobile را جستجو می کنیم تنها گروه معنایی که مشاهده می شود را نشان داده ایم.

$$\{car_n^1, auto_n^1, automobile_n^1, machine_n^4, motorcar_n^1\}$$

شماره بالا هر کلمه نشان دهنده این است در صورتی این کلمه را جستجو کنیم این مجموعه مترادف چندمین گروه مفهومی این کلمه است و شماره پایین نشان دهنده نقش کلمه در جمله است. مانند: اسم، فعل، پس همانطور که مشاهده کردید مجموعه مترادف یک مجموعه از مفهوم های کلمه است که کلمات هر دسته دارای معنا های تقریباً مشابهی هستند. بر اساس رابطه ای که در قسمت های قبلی درباره پیدا کردن مفهوم درست کلمه ها بر اساس نقشی را که دارا هستند توضیح دادیم، اکنون آن رابطه را برای یافتن مفهوم در وردنت به صورت زیر تغییر می دهیم.

$$Senses_{WN} : L \times POS \rightarrow \mathcal{P}^{SYNSETS} \quad (1)$$

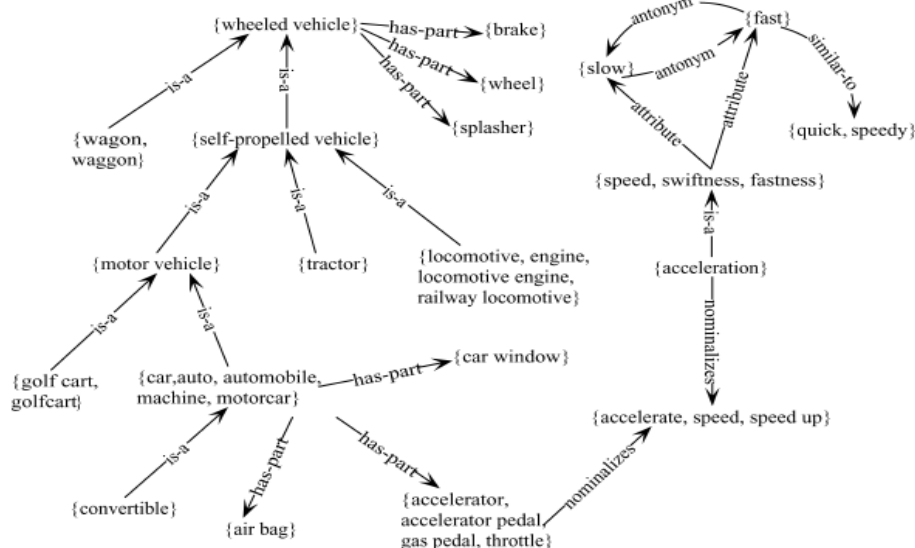
$$Senses_{WN}(car_n) = \left\{ \begin{array}{l} car_n^1, auto_n^1, automobile_n^1, machine_n^4, motorcar_n^1 \\ car_n^2, rail\ car_n^1, rail\ way\ car_n^1, rail\ road\ car_n^1 \\ cable\ car_n^1, car_n^3 \\ car_n^4, gondola_n^3 \\ car_n^3, elevator\ car_n^1 \end{array} \right\}$$

در صورتی که هر کلمه بدون ابهام باشد فقط باید یک مجموعه مترادف داشته باشد. به عنوان مثال  $car_n^1$  فقط با یک دسته معنایی مشخص ارتباط دارد.

$$\{car_n^1, auto_n^1, automobile_n^1, machine_n^4, motorcar_n^1\}$$

در شکل ۲ بخشی از شبکه معنایی وردنت که از مجموعه مترادف،  $car_n^1$  بدست آمده را نشان می دهد.





شکل ۲- ساختار داخلی وردنت همراه با رابطه ها (Navigli, ۲۰۰۹)

مثال (۲)

$$\{ \text{pop}_n^2, \text{soda}_n^2, \text{soda pop}_n^1, \text{soda water}_n^2, \text{tonic}_n^2 \}$$

عدهای بالای کلمات بیان کنندهی این است که در صورتی که کلمه ی مورد نظر را جستجو کنیم. آن مجموعه مترادف چندمین گروه دارای مفهوم کلمه ی مورد نظر است و عدهای پایین بیان کنندهی نقش آن کلمه است و  $n$  آخر کلمات مخفف اسم است و وردنت معنی کلمات را فراهم می کند. برای هر مجموعه مترادف یک مجموعه اطلاعات وجود دارد. یعنی رابطه هایی وجود دارد که اطلاعاتی را برای تکمیل کردن این داده در اختیار کاربرها قرار می دهد. تعدادی از این رابطه ها به شرح زیر است:

#### • Gloss

متن هایی است که برای بیان کردن بعضی از کارکردها و مثال هایی از این مجموعه مترادف می باشد.

#### • ارتباطات لغوی و معنایی

این قسمت برای برقرار کردن ارتباط بین مفاهیم و مجموعه مترادف ها است. تعدادی از ارتباطات معنایی و واژه ای که در وردنت استفاده شده اند به شرح زیر است:

**Antonymy:**  $X$  با  $Y$  دارای رابطه است اگر و فقط اگر این دو نسبت به یکدیگر دارای مفهوم متضاد باشند.

**Pertainymy:**  $X$  یک صفت است که می تواند بیان کند به صورتی به  $Y$  وابسته است.  $Y$  می تواند اسم یا یک صفت دیگر باشد.

**Nominalization:** این رابطه تبدیل اسم  $X$  به فعل  $Y$  را نمایش می دهد.



Hypernymy: این مورد، به نام های is-a یا kind-of نیز شناخته می شود. در این رابطه می گوئیم Y دارای رابطه Hypernymy با X است اگر و فقط اگر X یک قسمتی از Y باشد.

hyponymy و troponymy: این رابطه ها بر عکس رابطه Hypernymy مورد استفاده قرار می گیرند.

Meronymy: این رابطه به نام part of شناخته می شود. Y دارای رابطه Meronymy با X است اگر و فقط اگر Y یک بخشی از X باشد.

Holonymy: X دارای رابطه Holonymy با Y است اگر و فقط اگر X یک بخشی از Y باشد. این رابطه برعکس رابطه meronymy است.

Entailment: فعل Y برای اجرا شدن مستلزم اجرا شدن فعل X است. اگر می خواهید X را انجام دهید باید Y را انجام دهید.

Similarity: صفت X شبیه صفت Y است.

Attribute: اسم X یک خصوصیت است که صفت Y برای بیان کردن ارزشی برای اسم مورد استفاده قرار می گیرد.

See also: این یک ارتباط برای برقراری ارتباطات بین صفت ها مورد استفاده قرار می گیرد.

## ۲-۴-۱ توسعه ی وردنت

رویکردی که مطرح است متشکل از دو فاز اصلی است:

- نگاشت است. که به طور خود کار بین ویکیپدیا و وردنت انتشار می یابد.
- روابط اتصال صفحات ویکیپدیا است که به وردنت متصل می شود.

## ۲-۴-۲ نگاشت ویکیپدیا به وردنت

(۲)

$$\mu: Senses_{Wiki} \rightarrow Senses_{WN},$$

طوری که، برای هر Wikipage :

$$w \in Senses_{Wiki}:$$

$$\mu(w) = \begin{cases} s \in Senses_{WN}(w) & \text{اگر یک پیوند را بتوان انتشار داد} \\ \epsilon & \text{در غیر این صورت} \end{cases}$$

الگوریتم نگاشت: به منظور پیوند هر صفحه و ویکیپدیا به یک مفهوم وردنت، طبق مقاله (Ponzetto and Navigli, ۲۰۱۰) یک الگوریتم جدید توسعه داده شده است که شبه کد آن در الگوریتم ۱ نشان داده شده است و در مراحل زیر انجام می شود:

```

Input:  $Senses_{wiki}, Senses_{WN}$ 
Output: a mapping  $\mu: Senses_{wiki} \rightarrow Senses_{WN}$ 

1: for each  $w \in Senses_{wiki}$ 
2:    $\mu(w) := \epsilon$ 
3: for each  $w \in Senses_{wiki}$ 
4:   if  $|Senses_{wiki}(w)| = |Senses_{WN}(w)| = 1$  then
5:      $\mu(w) := w_n^1$ 
6: for each  $w \in Senses_{wiki}$ 
7:   if  $\mu(w) = \epsilon$  then
8:     for each  $d \in Senses_{wiki}$  s.t.  $d$  redirects to  $w$ 
9:       if  $\mu(d) \neq \epsilon$  and  $\mu(d)$  is in a synset of  $w$  then
10:         $\mu(w) := \text{sense of } w \text{ in synset of } \mu(d)$ ; break
11: for each  $w \in Senses_{wiki}$ 
12:   if  $\mu(w) = \epsilon$  then
13:     if no tie occurs then
14:        $\mu(w) := \underset{s \in Senses_{WN}(w)}{\operatorname{argmax}} p(s|w)$ 
15: return  $\mu$ 

```

### الگوریتم ۱- الگوریتم نگاشت (Ponzetto and Navigli, ۲۰۱۰)

در این الگوریتم، بهترین مفهوم  $S$ ، به وسیله محاسبه‌ی بخش‌هایی از زمینه‌ی ابهام زدایی از  $S$  و  $w$  و نرمال سازی توسط امتیازات خلاصه بر روی همه‌ی مفاهیم از  $w$  در ویکیپدیا و وردنت، بدست می‌آید. نتایج نشان داده است که بهبود روش‌های مطرح شده بر پایه‌ی حاشیه نویسی بزرگ و دارای عملکرد بالاتر، توسط استفاده از اطلاعات ابهام زدایی بیشتر، می‌تواند بدست آید و با استفاده از زمینه ابهام زدایی غنی کمک می‌کند تا مناسبترین مفهوم وردنت را برای یک صفحه ویکیپدیا انتخاب کنیم.

### ۳- منابع بدون ساختار

منابع بدون ساختار به گروه‌های زیر تقسیم می‌شود:

#### ۳-۱ Corpora

یک مجموعه از متون است که برای مدل‌های یادگیری استفاده می‌شود و شامل ۲ مدل دیگر می‌باشد که برای ابهام زدایی مفهوم کلمات با روش‌های بانظر و روش‌های بدون ناظر کاربرد دارند.

#### ۳-۱-۱ Raw corpora

مجموعه براون<sup>۵</sup> در سال ۱۹۶۷ ارائه شد و دارای یک میلیون کلمه در قالب مجموعه متن‌های متوازن است که در سال ۱۹۶۱ در آمریکا منتشر شده است. در حال حاضر، این مجموعه بالغ بر ۳۰ میلیون کلمه را شامل می‌شود.

#### ۳-۱-۲ Sense-Annotated Corpora

منظور از Sense Annotated یعنی کلیه مفاهیمی که از یک کلمه برداشت می‌شود را دارا است. این منبع برای بهبود در انجام ابهام زدایی مفهوم کلمات مورد استفاده قرار می‌گیرد. یکی از معروفترین منابعی که از این ساختار استفاده می‌کند

<sup>۵</sup>Brown

Semcor است.

### SemCor ۱-۲-۱-۳

Semcor زیر مجموعه ایی از مجموعه براون است که شامل کلماتی است که به طور دستی دارای ضمیمه های نقش کلمات در جمله، ریشه کلمه و مفهوم کلمه از دیدگاه وردنت است. Semcor دارای ۳۵۲ متن است که در ۱۸۶ تا از این متن ها، نقش کلمات مشخص شده است ولی در ۱۶۶ تا باقیمانده فقط کلمات فعل مشخص شده اند و مفهومی را که براساس نقش آن در جمله دارند نیز مشخص گردیده است. علاوه بر این، این مجموعه دارای ۲۳۴۰۰۰ تفسیر درباره ی مفهوم کلمات است. بنابراین یک مجموعه بزرگ از برچسب های مفهومی است که می توان آن را برای آموزش دادن روش هایی که برپایه باناظر برای ابهام زدایی هستند استفاده کرد. در شکل ۳ یک قسمت کوچک از متن هایی که در این مجموعه وجود دارد و برچسب گذاری شده است نشان داده شده است. همان طور که مشاهده می کنید همه ی کلمات بر اساس وردنت برچسب گذاری شده اند. به عنوان مثال کلمه ی word در جمله اول براساس وردنت دارای مفهوم "عبارت مختصر" و در جمله دوم دارای مفهوم "یک بخشی از زبان است که مردم محلی امکان تشخیص آن را دارند." است (حسامی و محمودی، ۱۳۸۹).

As of Sunday<sup>1</sup> night<sup>1</sup> there was<sup>4</sup> no word<sup>2</sup> of a resolution<sup>1</sup> being offered<sup>2</sup> there<sup>1</sup> to rescind<sup>1</sup> the action<sup>1</sup>. Pelham pointed out<sup>1</sup> that Georgia<sup>1</sup> voters<sup>1</sup> last<sup>1</sup> November<sup>1</sup> rejected<sup>2</sup> a constitutional<sup>1</sup> amendment<sup>1</sup> to allow<sup>2</sup> legislators<sup>1</sup> to vote<sup>1</sup> on pay<sup>1</sup> raises<sup>1</sup> for future<sup>1</sup> Legislature<sup>1</sup> sessions<sup>2</sup>.

### شکل ۳- ساختار داخلی Semcor همراه با برچسب گذاری کلمات (Navigli, ۲۰۰۹)

Semcor براساس نسخه وردنت ۱,۵ است ولی نسخه جاری از نسخه ۲,۱ استفاده کرده است. براساس semcor یک مجموعه دو زبانی به وسیله پیانتا<sup>۱</sup> و همکارانش به وجود آمده است که نام آن را Multisemcor نامیده اند. این مجموعه دارای دو مجموعه موازی به زبان های ایتالیایی و انگلیسی است. این مجموعه دارای برچسب های نقش کلمه، ریشه آنها و همچنین مفهوم آنها در ایتالیایی و انگلیسی است که برای مفاهیم از نسخه انگلیسی و ایتالیایی وردنت استفاده کرده اند. برچسب های مفهومی کلمات اصلی متن semcor به ایتالیایی ترجمه می شوند.

### ۴- پیش پردازش های ابهام زدایی مفهوم کلمات

پیش پردازش های ابهام زدایی مفهوم کلمات به شرح زیر است:

#### ۴-۱) تفکیک کردن

تفکیک کردن به این معنی است که یک مجموعه را به توکن هایی که سازنده این مجموعه است تقسیم بندی کنیم.

#### ۴-۲) برچسب زنی نقش کلمات

<sup>۱</sup>Pianta

در اینجا باید نقشی را که هر کلمه در این عبارت بازی می کند را تعیین کنیم. هر کلمه که در جمله قرار می گیرد دارای نقشی مانند اسم، فعل، صفت یا غیره است. همانطور که می دانیم کلمه در هر نقشی که قرار می گیرد مفهوم خاصی را پیدا می کند.

#### ۴-۳ ریشه یابی

در این گام ریشه هر کلمه را بدست می آوریم. این کار به علت آنکه هر کلمه در هر موقعیت زمانی از لحاظ گرامر قرار می گیرد به آن پیشوند و پسوندهایی می چسبد و حتی در مواقعی نیز موجب تغییر در ریشه کلمه می شود، اما همه ی این تغییرات باز هم آن کلمه دارای مفهوم ریشه خودش است. برای اینکه کار شناسایی مفهوم کلمه آسان تر شود باید تمام کلمه هایی که در جمله قرار دارند را ریشه یابی نماییم. برای اینکار روش های مختلفی مانند ریشه یاب پورتر برای زبان انگلیسی و ریشه یاب کاظم تقوی برای زبان فارسی ارائه شده اند.

#### ۴-۴ قطعه بندی

بر اساس قواعد زبان شناسی، زمانی که دو یا چند کلمه در کنار یکدیگر قرار می گیرند موجب ایجاد یک عبارت قاعده ایی خاص می شوند. به عنوان مثال عبارت اسمی یا عبارت فعلی. این کار یکی از کارها در پردازش زبان های طبیعی است که در آن برای بدست آوردن نوع عبارت مورد استفاده قرار می گیرد.

#### ۴-۵ تجزیه کردن

تعیین ساختار نحوی جمله که به طور متداول به کمک رسم درخت پارسر است. که در آن درختی را بر اساس داده هایی که داریم می سازیم تا ببینیم که آیا ساختار این جمله درست است یا خیر.

#### ۵- انتخاب متد کلاس بندی

آخرین مرحله ایی که در ابهام زدایی مفهوم کلمات مطرح است انتخاب کردن یکی از روش های کلاس بندی است. بیشترین روش هایی که در این مرحله استفاده می شود روش های یادگیری ماشین است. همچنین از روش های شناسایی الگو نیز می توان استفاده کرد ولی استقبال چندانی از آنها نشده است. به طور کلی دو دسته بندی برای روش های ابهام زدایی مفهوم کلمات وجود دارد روش های باناظر که در این روش ها ابتدا از یک مجموعه برای آموزش استفاده می کنیم و سپس بر اساس اطلاعات بدست آمده کار دسته بندی های بعدی را انجام می دهیم. روش هایی که از بدون ناظر استفاده می کنند آموزشی در کار نیست، بلکه در این روش ها از یک منبع دانشی استفاده اطلاعات مناسب را از داخل منبع دانش استخراج می کنند. همچنین پایگاه های دانشی که مورد استفاده قرار می گیرند را هم می توان به دو دسته تقسیم بندی کرد، پایگاه دانش قوی و ضعیف. در پایگاه های دانش قوی از ساختارهایی مانند آنتولوژی استفاده می شود که در آنها تمام رابطه های بین واژگان در نظر گرفته شده است. از این نوع پایگاه های دانش می توان وردنت را نام برد. اما پایگاه های دانش ضعیف مانند براون که مجموعه از تمامی کلمات موجود است که همه آنها را برچسب گذاری کرده اند و ارتباطات بین کلمه ها مشخص نشده است. این ساختار دارای یک

قسمت به نام حذف کلمات اضافه است. این قسمت مربوط به کلماتی در جمله است که در اصل دارای بار معنایی خاصی نیست و فقط برای اینکه خواننده راحتتر بتواند متن را درک نماید مورد استفاده قرار می گیرند. به عنوان مثال بعضی از این کلمات در زبان انگلیسی عبارتند از the, a, it, that و غیره. در ابهام زدایی مفهوم کلمات این کلمات در نظر گرفته نمی شوند تا بتوان سریع تر و راحتتر مفهوم کلمات اصلی را بدست آورد.

## ۶- معرفی چهار دسته از روش های یادگیری در ابهام زدایی مفهوم کلمات

در این بخش، چهار روش یادگیری در ابهام زدایی مفهوم کلمات مورد بررسی قرار می گیرد که شامل: روش با ناظر<sup>۷</sup>، روش بدون ناظر<sup>۸</sup>، روش نیمه ناظر<sup>۹</sup>، روش مبتنی بر پایگاه دانش و فرهنگ لغت<sup>۱۰</sup> است.

### ۶-۱- روش با ناظر

در روش های با ناظر از یک مجموعه داده آموزشی برای آموزش استفاده شده و سپس براساس این داده آموزشی تصمیم گیری های مناسب صورت می پذیرد. روش های با ناظر بر روی مسئله ابهام زدایی مفهوم کلمات، دارای کارایی بالایی است و نتایج بسیار خوبی را نیز بدست آورده است ولی دارای مشکلاتی است. مانند: تهیه کردن یک مجموعه داده آموزشی که بتواند کلیه یا اکثریت مفاهیم کلمه ها را تحت پوشش خودش داشته باشد. این روش، به طور عمده زمینه ای از واژه ها را برای ابهام زدایی انتخاب می کند. روش با ناظر شامل مراحل تحت نظارت و آموزش و مرحله ی آزمایش است. روش های با ناظر گران هستند و به دلیل کمبود اطلاعات، شکننده و انعطاف ناپذیرند و برای حاشیه نویسی و شرح اطلاعات دشوار است. ابهام زدایی مفهوم کلمات مبتنی بر یادگیری روش با ناظر بهترین عملکرد را در SenseEval و کارگاه های SemEval دارند (Zhong et al, ۲۰۰۸). الگوریتم های یادگیری با ناظر شامل: روش کلاسه بندی ماشین بردار پشتیبانی، درخت تصمیم گیری، شبکه های عصبی، لیست تصمیم گیری، بیزین، k نزدیک ترین همسایه و بوستینگ است.

### ۶-۲- روش بدون ناظر

روش های مبتنی بر بدون ناظر با روش های با ناظر و مبتنی بر پایگاه دانش بسیار متفاوت هستند. به عبارتی دیگر در روش های بدون ناظر هدف تعیین کردن دسته هر کدام از مفاهیم است. ولی روش های دیگر برای برچسب گذاری هر کدام از کلمه ها استفاده می شود. به هر حال هر دو مورد تفکیک کردن مفاهیم و برچسب گذاری از چالش های موجود در ابهام زدایی مفاهیم است. در سال های اخیر، روش بدون ناظر، در استفاده از منابع دانش گسترده با دامنه ی خاصی از اطلاعات، موفق آمیز بوده

<sup>۷</sup> Supervisor

<sup>۸</sup> Unsupervised

<sup>۹</sup> Semi-supervised

<sup>۱۰</sup> Dictionary and knowledge based methods

است که می توان به فعالیت جن<sup>۱۱</sup> و همکارانش در سال ۲۰۰۹ اشاره داشت (Raviv and Markovitch, ۲۰۱۲). در سال ۲۰۰۱، بریل<sup>۱۲</sup> و همکارانش روش های بدون ناظر، نیمه ناظر و یا یادگیری فعال را با استفاده از مجموعه داده های بزرگ، هنگامی که برچسب ها گران قیمت بود پیشنهاد کردند (Yuret and Yatbaz, ۲۰۱۰). روش بدون ناظر بسیار قوی و قابل حمل<sup>۱۳</sup> و در واقع به طور دستی است و به منابعی مانند مفاهیم سلسله مراتبی<sup>۱۴</sup> و واژه نامه ها و منابع دانش دستی نیاز ندارد (Tomar et al, ۲۰۱۳). الگوریتم های ابهام زدایی مفهوم کلمات بدون ناظر به دو دسته ی کلی تقسیم می شوند:

- ابهام زدایی مفهوم کلمات مبتنی بر Token<sup>۱۵</sup>: کسانی که عمل ابهام زدایی مبتنی بر Token را به وسیله کاوش شباهت ها و یا ارتباط بین یک کلمه ی مبهم و متن آن را می یابند.
- ابهام زدایی مفهوم کلمات مبتنی بر نوع<sup>۱۶</sup>: آنهایی که عمل ابهام زدایی مفهوم کلمات را بر اساس نوع، به سادگی با اختصاص تمام موارد و مثال های رایج برای یک کلمه ی مبهم انجام می دهند (Brody et al, ۲۰۰۶).

### ۶-۳- روش نیمه ناظر

یادگیری نیمه ناظر بین یادگیری بدون ناظر و یادگیری باناظر قرار گرفته و در اینجا بسیار مناسب است. زیرا طبقه بندی بهتری را با استفاده از تعداد زیادی از داده های بدون برچسب و تعداد محدودی داده ی برچسب دار ایجاد می کند. در واقع یادگیری نیمه ناظر، از هر دوی داده های برچسب دار و بدون برچسب برای آموزش استفاده می کند. در حقیقت تعداد داده های بدون برچسب در این نوع یادگیری نسبت به داده های برچسب دار بسیار بیشتر است. تحقیقات در زمینه ی یادگیری ماشین وجود دارد که نشان داده اند زمانی که داده های بدون برچسب به همراه تعداد اندک داده برچسب دار به کار می روند، می توانند بهبود چشمگیری در دقت یادگیری ایجاد نمایند (عارفیان و افتخاری، ۱۳۹۲). بسیاری از الگوریتم های نیمه ناظر مختلف برای وظایف پردازش زبان طبیعی استفاده شده اند (Søgaard, ۲۰۱۱). نام یادگیری نیمه ناظر، حاکی از این واقعیت است که، داده بین یادگیری باناظر (با داده های آموزشی برچسب دار) و یادگیری بدون ناظر (بدون هیچ گونه داده های آموزشی با برچسب) مورد استفاده قرار می گیرد. بسیاری از محققان یادگیری ماشین یافتند که داده های بدون برچسب، زمانی که در ارتباط با مقدار کمی از داده های برچسب دار استفاده می شود، می تواند بهبود قابل توجهی در دقت یادگیری ایجاد کند. کسب داده های برچسب

<sup>۱۱</sup>chen

<sup>۱۲</sup>Brill

<sup>۱۳</sup>protable

<sup>۱۴</sup>hierarchies concept

<sup>۱۵</sup>Token\_based WSD

<sup>۱۶</sup>Type\_based WSD

دار برای مشکل یادگیری، اغلب به یک عامل انسانی ماهر و فرآیند گران قیمت و زمانبر نیاز دارد. در حالی که دستیابی به داده های بدون برچسب نسبتاً ارزان است و وقت کمتری صرف می کند (Sati, ۲۰۱۳).

## ۶-۴- روش مبتنی بر پایگاه دانش و فرهنگ لغت

این روش با استفاده از دانش واژگانی<sup>۱۷</sup> مانند واژه نامه ها و اصطلاحنامه و فرض اینکه دانش را می توان از تعاریفی از کلمات استخراج کرد. در واقع در این روش، اگر اطلاعاتی درباره ی طبقه بندی معنایی یک کلمه وجود نداشته باشد، می توان از مشخصات عمومی معنای یک کلمه در فرهنگ لغت استفاده کرد. به عنوان مثال الگوریتم لسک که در این روش اول مفاهیم کلمه های مورد نظر را مقایسه می کند و برای محاسبه ی همپوشانی مفاهیم، کلمات مورد نظر استفاده می شود و با بیشترین همپوشانی<sup>۱۸</sup> کلمه در تعریف فرهنگ لغت خود همراه است. روش های مبتنی بر دانش، به ابهام زدایی مفهوم کلمات با مقایسه ی متن خود با اطلاعات از یک منبع واژگانی از پیش تعریف شده مانند وردنت می پردازد (Cabrera et al, ۲۰۰۹).

## ۷- نتیجه گیری

یکی از اولین مشکلاتی که در هر سیستم پردازش زبان طبیعی با آن برخورد می کنیم، مسئله ابهام معنایی و ساختاری کلمات است. یافتن مفهوم صحیح کلمات در یک جمله برای یک ماشین بسیار دشوار است. زیرا ماشین باید عبارت های انسان را که بدون ساختار هستند، به عبارت های با ساختار تبدیل نماید تا بتواند مفهوم صحیح کلمات را تشخیص دهد که به این تشخیص، ابهام زدایی مفهوم کلمات گفته می شود. روال اصلی هر سیستم ابهام زدایی دارای دو قسمت است: اول، یک پایگاه دانش و دوم، یک الگوریتمی که بتواند معنی صحیح را از پایگاه دانش استخراج نماید.

منابع به دو دسته منابع دارای ساختار و منابع بدون ساختار تقسیم بندی می شوند و برای کسب مفهوم صحیح کلمات باید بر روی آن پیش پردازش هایی انجام شود تا بهترین مفهوم بدست آید که در این مقاله بطور جامع به بررسی انواع منابع دارای ساختار و منابع بدون ساختار و پیش پردازش های ابهام زدایی مفهوم کلمات در پردازش زبان طبیعی پرداختیم و بر اساس بررسی های انجام شده، استفاده از وردنت پیشنهاد می شود که یک منبع ضروری برای ابهام زدایی مفهوم کلمات است و یک منبع دارای ساختار می باشد.

## منابع

<sup>۱۷</sup> Lexical knowledge

<sup>۱۸</sup> Overlap



دکتر مهرنوش شمس فرد، دکتر احمد عبدالله زاده بارفروش. "استخراج دانش مفهومی از متن با استفاده از الگوهای زبانی و معنایی". تازه های علوم شناختی، سال ۴، شماره ۱، ۱۳۸۱.  
مرتضی نوروزی، محسن طاهریان. "وب معنایی". مجری طرح: سازمان فناوری اطلاعات ایران، مدیر طرح: محمد رضا فروزنده دوست. چاپ اول ۱۳۹۰.

احسان حسامی، دکتر محمودی. "ابهام زدایی مفهوم کلمات". ارائه شده برای: سمینار کارشناسی ارشد مهندسی نرم افزار، دانشگاه آزاد اسلامی واحد قزوین، دانشکده برق، رایانه و فن آوری اطلاعات. سال تحصیلی ۸۹-۹۰.

فاطمه عارفیان، مهدی افتخاری. "روش جدید K نزدیکترین همسایه فازی و ناهموار برای طبقه بندی نیمه نظارتی". مهندسی کامپیوتر و توسعه پایدار با محوریت شبکه های کامپیوتری، مدل سازی و امنیت سیستم ها. برگزار کننده: موسسه آموزش عالی خاوران مشهد، ۲۸ آذرماه ۱۳۹۲.

E. Kilgarriff, "Introduction to the special issue on evaluating word sense disambiguation systems", Journal of Natural Language Engineering, ۸(۴): ۲۷۹-۲۹, ۲۰۰۲.

Ch.Ho, M.Murad, R.Kadir, Sh.Doraisamy. "Word Sense Disambiguation-based Sentence Similarity". Coling ۲۰۱۰: Poster Volume, pages ۴۱۸-۴۲۶, Beijing, August ۲۰۱۰.

T.Miller, N.Erbs, H.Zorn, T.Zesch, I.Gurevych. "DKPro WSD – A Generalized UIMA-based Framework for Word Sense Disambiguation". Proceedings of the ۵۱st Annual Meeting of the Association for Computational Linguistics, pages ۳۷-۴۲, Sofia, Bulgaria, August ۴-۹ ۲۰۱۳. ©۲۰۱۳ Association for Computational Linguistics.

T.Pedersen, R.Bruce. "Knowledge Lean Word-Sense Disambiguation". AAAI-۹۸ Proceedings. Copyright © ۱۹۹۸, AAAI (www.aaai.org). All rights reserved.

Z.Zhong, H.Tou, Y.Chan. "Word Sense Disambiguation Using OntoNotes: An Empirical Study". Proceedings of the ۲۰۰۸ Conference on Empirical Methods in Natural Language Processing, pages ۱۰۰۲-۱۰۱۰, Honolulu, October ۲۰۰۸. ©۲۰۰۸ Association for Computational Linguistics.

R.Navigli, "Word Sense Disambiguation: A Survey", ACM Computing Surveys, Vol. ۴۱, No. ۲, Article ۱۰, Publication date: February ۲۰۰۹.

J.Chen. "Adaptive Word Sense Disambiguation Using Lexical Knowledge in a Machine-readable Dictionary". Computational Linguistics and Chinese Language Processing, Vol. ۵, No. ۲, August ۲۰۰۰, pages ۲-۵. © Computational Linguistics Society of R.O.C.

L.Tan. "Examining Crosslingual Word Sense Disambiguation". A thesis submitted to the Nanyang Technological University in partial fulfillment of the requirement for the

- degree of Masters of Arts, CHAPTER ۴. KNOWLEDGE SOURCES FOR WSD, pages ۲۸-۲۹. School of Humanities and Social Sciences. ۲۰۱۳.
- S.Patwardhan, S.Banerjee, and T.Pedersen. "Using Measures of Semantic Relatedness for Word Sense Disambiguation". A. Gelbukh (Ed.): CICLing ۲۰۰۳, LNCS ۲۵۸۸, pp. ۲۴۱-۲۵۷, ۲۰۰۳. ©Springer-Verlag Berlin Heidelberg ۲۰۰۳.
- T.Tanaka, F.Bond, T.Baldwin, S.Fujita, Ch.Hashimoto. "Word Sense Disambiguation Incorporating Lexical and Structural Semantic Information". Proceedings of the ۲۰۰۷ Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, pp. ۴۷۷-۴۸۵, Prague, June ۲۰۰۷. ©۲۰۰۷ Association for Computational Linguistics.
- A.Philpot, E.Hovy, P.Pantel, "The Omega Ontology", In Proceedings of the IJCNLP Workshop on Ontologies and Lexical Resources (OntoLex, Jeju Island, South Korea), ۲۰۰۵.
- S.Ponzetto, R.Navigli. "Knowledge-rich Word Sense Disambiguation Rivaling Supervised Systems". Proceedings of the ۴۸th Annual Meeting of the Association for Computational Linguistics, pages ۱۵۲۲-۱۵۳۱, Uppsala, Sweden, ۱۱-۱۶ July ۲۰۱۰.
- D.Yuret, M.Yatbaz. "The Noisy Channel Model for Unsupervised Word Sense Disambiguation". ©۲۰۱۰ Association for Computational Linguistics. Volume ۳۶, Number ۱. pages ۱۱۲.
- A.Raviv, Sh.Markovitch. "Concept-Based Approach to Word-Sense Disambiguation". Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence, Copyright ©۲۰۱۲, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.
- G.Tomar, M.Singh, Sh.Rai, A.Kumar, R.Sanyal, S.Sanyal. "Probabilistic Latent Semantic Analysis for Unsupervised Word Sense Disambiguation". IJCSI International Journal of Computer Science Issues, Vol. ۱۰, Issue ۵, No ۲, September ۲۰۱۳ ISSN (Print): ۱۶۹۴-۰۸۱۴, ISSN (Online): ۱۶۹۴-۰۷۸۴ [www.IJCSI.org](http://www.IJCSI.org). Copyright (c) ۲۰۱۳ International Journal of Computer Science Issues. All Rights Reserved.
- S.Brody, R.Navigli, M.Lapata. "Ensemble Methods for Unsupervised WSD". Proceedings of the ۲۱st International Conference on Computational Linguistics and ۴۴th Annual

- Meeting of the ACL, pages ۹۷-۱۰۴, Sydney, July ۲۰۰۶. ©۲۰۰۶ Association for Computational Linguistics.
- Ms. Ankita Sati. "Review: Semi-Supervised Learning Methods for Word Sense Disambiguation". IOSR Journal of Computer Engineering (IOSR-JCE) e-ISSN: ۲۲۷۸-۰۶۶۱, p- ISSN: ۲۲۷۸-۸۷۲۷ Volume ۱۲, Issue ۴ (Jul. - Aug. ۲۰۱۳), PP ۶۳-۶۸. [www.iosrjournals.org](http://www.iosrjournals.org).
- A.Søgaard. "Semisupervised condensed nearest neighbor for part-of-speech tagging". Proceedings of the ۴۹th Annual Meeting of the Association for Computational Linguistics: shortpapers, pages ۴۸-۵۲, Portland, Oregon, June ۱۹-۲۴, ۲۰۱۱. ©۲۰۱۱ Association for Computational Linguistics.
- R.Cabrera, P.Rosso, M.Gómez, L.Pineda, and D.Avendaño. "Semi-supervised Word Sense Disambiguation Using the Web as Corpus". A. Gelbukh (Ed.): CICLing ۲۰۰۹, LNCS ۵۴۴۹, pp. ۲۵۶-۲۶۵, ۲۰۰۹. © Springer-Verlag Berlin Heidelberg ۲۰۰۹.