



بازشناسی رابطه تقابل در زبان فارسی در حالت حذف نقش‌نمای گفتمان به کمک ماشین بردار پشتیبان

حبیب خدادادی^(۱) - سعید راحتی قوچانی^(۲) - اعظم استاجی^(۳)

(۱) گروه کامپیوتر - دانشگاه آزاد اسلامی واحد مشهد

habibekhodadi@gmail.com

(۲) گروه برق - دانشگاه آزاد اسلامی واحد مشهد

Rahati@Mshdiau.ac.ir

(۳) گروه زبان‌شناسی - دانشگاه فردوسی مشهد

estagi@ferdowsi.um.ac.ir

خلاصه: شناسایی رابطه تقابل در یک گفتمان می‌تواند به کمک نقش‌نماهای خاص رابطه تقابل مانند "اما" و "ولی" انجام شود؛ اما این نقش‌نماها در مواردی حذف می‌شوند و شناسایی رابطه را با مشکل مواجه می‌کنند. به همین علت از ویژگی‌هایی شامل زمان فعل‌ها، استفاده از جفت کلمات، وجود اعداد در دو قسمت متنی اطراف نقش‌نما و وجود فعل‌ها و کلمات منفی در یک طرف نقش‌نما و عدم وجود آن در طرف دیگر به منظور شناسایی این رابطه استفاده شد. در این مقاله از مجموعه داده پژوهش‌شده هوشمند علائم استفاده شده و ۵۰۰۰ نمونه رابطه تقابل و ۵۰۰۰ نمونه سایر روابط گردآوری شده، بردار ویژگی را به ازای هر نمونه شکل داده و دسته‌بند *SVM* با این ویژگی‌ها آموزش داده شد. میزان صحت دسته‌بند ۶۱.۰۲ در بهترین حالت می‌باشد.

کلمات کلیدی: گفتمان، نقش‌نمای گفتمان، شناسایی رابطه تقابل، ماشین بردار پشتیبان.

۱ - مقدمه

"پس"، "بنابراین"، و ... را نام برد. در صورت وجود داشتن نقش نماها، از آنها می‌توان به عنوان عامل شناسایی نوع رابطه استفاده کرد ولی در مواردی این عناصر وجود ندارند.

در تحقیقات گسترده‌ای که به خصوص در زبان انگلیسی صورت گرفته است، با استخراج ویژگی‌هایی سعی در شناسایی این روابط در سطح گفتمان شده است؛ اما در زبان فارسی در این مورد کاری انجام نشده است.

در [۳] با جمع‌آوری نمونه‌های زیادی از ۴ رابطه گفتمانی و سپس حذف نقش‌نماها گفتمان و با این فرض که خود کلمات موجود در رابطه می‌توانند نوع رابطه را تعیین کنند، یک مدل احتمالاتی از جفت کلمات دو طرف رابطه را بدست آورده (هر جفت کلمه شامل یک کلمه قبل و یک کلمه بعد از نقش‌نما می‌باشد) و سپس از این مدل و با داشتن کلمات دو طرف هر رابطه گفتمانی و بدون استفاده از نقش‌نما، احتمال تعلق به هر کدام از رابطه‌ها را بدست آورده و بیشترین احتمال نوع رابطه را تعیین می‌کند.

در [۴] که بر روی زبان ژاپنی کار شده است، علاوه بر استفاده از جفت

گفتمان (*Discourse*) به هر قطعه یا پاره‌ای از زبان گفته می‌شود که به قصد برقراری ارتباط به کار برده شده باشد [۱] و به عنوان واحد بالاتر از جمله در زبان‌شناسی شناخته می‌شود. پردازش کامپیوتری گفتمان یکی از شاخه‌های نسبتاً جدید حوزه پردازش متن می‌باشد که در زبان فارسی به آن چندان پرداخته نشده است. یکی از زمینه‌های مطالعه در پردازش گفتمان، شناسایی و برچسب‌گذاری روابط بین واحدهای متنی (بند، جمله یا پاراگراف) در سطح یک گفتمان است. بعضی از این رابطه‌ها شامل علی، افزایشی، استمرار، تقابل، زمانی و ... می‌باشد.

نقش‌نماهای گفتمان (*Discourse connectives*) عناصری در سطح گفتمان هستند که کارکرد آنها نشان دادن رابطه‌ای است که بین یک پاره گفتار و پاره گفتار قبلی وجود دارد [۱]. از جمله این عناصر در زبان انگلیسی می‌توان "And"، "Or"، "But"، "Well"، "Because" و ... و در زبان فارسی "اما"، "ولی"، "خب"، "حالا"،

نقش نمای *since* هم برای رابطه زمانی (*temporal*) و هم برای رابطه علی (*causal*) به کار می رود [۹].

۳- در خیلی از موارد نقش نماها حذف می شوند (مثال ۴) [۶].
مثال ۲:

(۲a) *Selling picked up as previous buyers bailed out of their positions and aggressive short sellers- anticipating further declines-moved in.*

(۲b) *my favorite colors are blue and green.*

مثال ۳:

(۳a) *there have been more than ۱۰۰ mergers and acquisitions Within the European paper industry since the most recent wave of friendly takeovers was completed in the U.S. in ۱۹۸۶.*

(۳b) *It was a far safer deal for lenders since NWA had a healthier cash flow and more collateral on hand.*

مثال ۴:

"A lot of investor confidence comes from the fact that they can speak to us," he says.

"To maintain that dialogue is absolutely crucial."

بنابراین در حالتی که هیچ یک از ۳ مشکل بالا برای یک نقش نما خاص اتفاق نیفتد، صرفاً با نقش نما می توان نوع رابطه را شناسائی کرد؛ اما در حالتی که یکی از مشکلات بالا برای یک نقش نمای خاص اتفاق می افتد باید از ابزارهای واژگانی، نحوی و... دیگری برای این منظور استفاده کرد.

یکی از روابط موجود در گفتمان در زبان فارسی رابطه تقابلی است. از نظر مناسبات معنایی در رابطه تقابلی در سطح گفتمان بین بخش اول و بخش دوم به نوعی نتیجه خلافی، عکس یا مغایرت وجود دارد. این رابطه در سطح گفتمان به وسیله نقش نما های گفتمان مخصوص رابطه تقابلی شناخته می شوند از قبیل اما، ولی، لکن، با این همه، با این حال، علی رغم و... این نقش نماها می توانند در مواردی حذف شوند که خواننده (شنونده) از روی معنا به نوع رابطه پی می برد. در مثال ۵ و ۶ دو نمونه از رابطه تقابلی از پیکره متنی بی جن خان (<http://ece.ut.ac.ir/dbrg/bijankhan>) نشان داده شده است که اولی دارای نقش نمای تقابلی "ولی" می باشد، اما دومی فاقد نقش نما می باشد.

مثال ۵:

البته تمام اقلام غذایی ذکرشده برای همه مبتلایان به میگرن اثر نامطلوب ندارند ولی در تعداد زیادی از بیماران موثر هستند.

مثال ۶:

این دوربین ابتدا برای تشخیص پرتوهای کیهانی طراحی شده بود، در حال حاضر با تغییراتی که در آن ایجاد شده است، مصارف پزشکی نیز پیدا کرده است.

بنابراین و با توجه به مشکل حذف نقش نمای تقابلی، از روش های

کلمات از اطلاعات الگوهای گروهی (*Phrasal Pattern*) نیز استفاده شده است. این الگوهای عبارتی قسمت هائی از جمله هستند که در بعضی از روابط ظاهر می شوند. مثلاً رابطه تضاد در زبان انگلیسی در مواردی به صورت "... did" And "... should have done ..." ظاهر می شود. با استفاده از این اطلاعات اقدام به شناسائی روابط شده است. در [۵] فقط به بررسی حالاتی که نقش نما وجود ندارد پرداخته شده است و بر روی PDTB، که مجموعه برچسب خورده گفتمانی در زبان انگلیسی است، نمونه های حذف نقش نما پیدا شده و ویژگی هائی مانند وجود قطب های مختلف کلمات، کلاس های فعلی، جفت کلمات و... استخراج شده و با کمک این ویژگی ها دسته بندی صورت گرفته است.

در [۶] نیز برای حالتی که نقش نما وجود ندارد و بر روی PDTB ویژگی هائی استخراج شده است و با کمک این ویژگی ها نوع رابطه تشخیص داده شده است. که از ویژگی هائی مانند روابط قبل و بعد رابطه فعلی، وجود وابستگی های رابطه ای، ویژگی های استخراج شده از درخت نحو و... در این کار استفاده شده است. در [۷] که از مجموعه داده PDTB۲ استفاده شده است، ضمن استفاده از ویژگی های مقالات قبلی ویژگی جدیدی معرفی شده است، که با استفاده از یک مدل زبانی ابتدا نقش نما پیش گوئی می شود (در حالت نبود نقش نما) و سپس از این نقش نمای پیش گوئی شده به عنوان یک ویژگی اضافی به منظور شناسائی نوع رابطه استفاده می شود.

۲- رابطه های گفتمانی و رابطه تقابلی

در سطح هر گفتمان روابطی بین واحدهای متنی وجود دارد که در حالت عادی نوع این روابط به کمک عناصری به نام نقش نماهای گفتمان تشخیص داده می شود. شناسائی رابطه های موجود در یک گفتمان کاربردهای گسترده ای دارد؛ از جمله این که شناسائی این روابط به توانائی تولید و درک گفتمان کمک می کند و در سیستم های پرسش و پاسخ و خلاصه ساز و... کاربرد دارند [۳]. مثال ۱ نمونه ای از یک رابطه گفتمانی در زبان انگلیسی می باشد که با وجود نقش نمای *But* نوع رابطه *Contrast* تشخیص داده می شود [۶].

مثال ۱:

In any case, the brokerage firms are clearly moving faster to create new ads than they did in the fall of ۱۹۸۷.

But it remains to be seen whether their ads will be any more effective.

تشخیص نوع رابطه به کمک نقش نما با سه مشکل مواجه است:

۱- بعضی از نقش نماها دارای کارکرد غیر نقش نمائی هم می باشند؛ مانند مثال ۲ قسمت ۲b که کلمه *and* نقش نمای گفتمان نیست ولی در ۲a نقش نمای گفتمان می باشد [۹].

۲- بعضی از نقش نماها بین دو یا چند رابطه مشترک هستند و تشخیص رابطه با ابهام مواجه می شود؛ مانند مثال ۳ که

دیگری به جز استفاده از نقش نما برای شناسایی این رابطه استفاده می کنیم تا در موقع حذف نیز مشکلی برای شناسایی وجود نداشته باشد.

۳- مشکلات بازشناسی روابط گفتمان در زبان فارسی

یکی از بزرگترین موانع پیش روی پردازش کامپیوتری گفتمان، نبود داده برچسب خورده خاص گفتمان می باشد. در زبان انگلیسی داده برچسب خورده وجود دارد. یکی از این داده ها، داده Penn Discourse Treebank (PDTB) [۸] می باشد؛ که ضمن برچسب گذاری هر رابطه، در حالت حذف نقش نما نیز، بهترین نقش نمائی که می تواند در آنجا قرار گیرد معرفی شده است. و برای هر نقش نما هم دو بخش قبل و بعد آن با Arg_1 و Arg_2 نشان داده شده است. مثال ۷ و ۸ دو نمونه از این داده ها را نشان می دهد. مثال ۷:

Arg₁: *In any case, the brokerage firms are clearly moving faster to create new ads than they did in the fall of ۱۹۸۷.*

Arg₂: *But it remains to be seen whether their ads will be any more effective.*

(Contrast)

مثال ۸:

Arg₁: *"A lot of investor confidence comes from the fact that they can speak to us," he says.*

Arg₂: *[so] "To maintain that dialogue is absolutely crucial."*

(Cause)

علاوه بر عدم وجود داده برچسب خورده، فقدان برخی از ابزارهای پردازش زبان طبیعی مانند wordnet و... نیز از مشکلات مهم می باشد. مشکلات طبیعی رسم الخط زبان فارسی از قبیل مشکل کلمات مرکب و راست به چپ نویسی و... نیز وجود دارد.

۴- استخراج ویژگی برای سیستم شناسایی رابطه تقابل در گفتمان های فارسی

در این مقاله ما به منظور شناسایی رابطه تقابل در زبان فارسی از ۴ ویژگی استفاده کرده ایم. در همه جا منظور از Arg_2 عبارت و جملاتی است که نقش نما به طور ساختاری در آن طرف قرار دارد و طرف دیگر به عنوان Arg_1 شناخته می شود. این ویژگی ها شامل:

۱- از آنجا که در عبارات تقابلی ممکن است کلمات و فعل های مثبت و منفی در دو طرف نقش نما ظاهر شوند، اولین ویژگی عبارت است از وجود فعل ها و کلمه های منفی در یک طرف نقش نما و عدم وجود آن در طرف دیگر.

۲- اگر هر دو Arg_1 و Arg_2 شامل عدد، درصد و یا مقدار عددی باشد، احتمالاً یک رابطه تقابل وجود دارد [۵]. بنابراین دومین ویژگی شامل وجود عدد، درصد و یا مقدار

عددی در هر دوی Arg_1 و Arg_2 می باشد.

۳- هم زمان بودن فعل های اصلی Arg_1 و Arg_2 . از آنجا که در بعضی از رابطه ها مانند رابطه زمانی (Temporal)، زمان فعل های Arg_1 و Arg_2 ممکن است با هم تفاوت داشته باشد؛ هم زمان بودن فعل قبل از نقش نما و آخرین فعل بعد از نقش نما به عنوان ویژگی در نظر گرفته شده است.

۴- وجود جفت کلمات پرکاربرد رابطه تقابلی.

اگر رابطه تقابل بین دو قسمت متنی برقرار باشد و W_1 و W_2 به تمام کلمات این دو قسمت اشاره کند، فرض ما این است که جفت کلمات $(w_i, w_j) \in W_1 \times W_2$ (حاصل ضرب برداری دو طرف رابطه) می تواند به ما وجود رابطه تقابل را نشان دهد. در این راستا و به منظور شناسایی این نوع جفت کلمات کارهای زیر انجام شد:

۱- از دو پیکره بی جن خان و پیکره ی متنی زبان فارسی که توسط پژوهشکده پردازش هوشمند علائم (http://www.rcisp.com) تهیه شده است، تمام نمونه های رابطه تقابل که دارای یکی از نقش نماهای "ولی"، "اما"، "لکن"، "باین وجود"، "باین حال"، "باوجود این" و "لیکن" هستند، استخراج شدند که کلمات قبل از آن به عنوان Arg_1 و کلمات بعد از آن به عنوان Arg_2 برچسب زده شد. مثال هایی از هر کدام از این نقش نماها را در جدول ۱ مشاهده می کنید.

۲- این نقش نماها و هر کلمه ای که نقشی بجز فعل، اسم، صفت یا قید دارد حذف شدند. همچنین فعل است هم که معنی خاصی را به همراه ندارد نیز حذف شد.

۳- هر کلمه از pos آن جدا شد.

۴- ضرب برداری کلمات باقی مانده در دو طرف رابطه انجام شد بدین معنی که هر کلمه در Arg_1 را با تک تک کلمات Arg_2 به صورت یک جفت مجزا در نظر گرفته و در تمام نمونه ها تعداد این جفت ها شمرده می شود.

۵- تمامی جفت کلماتی که زیر ۵۰ بار تکرار شده است را حذف و باقی جفت کلمات را ذخیره می کنیم. بنابراین وجود داشتن این جفت کلمات در هر رابطه را به عنوان یک ویژگی در نظر گرفته ایم.

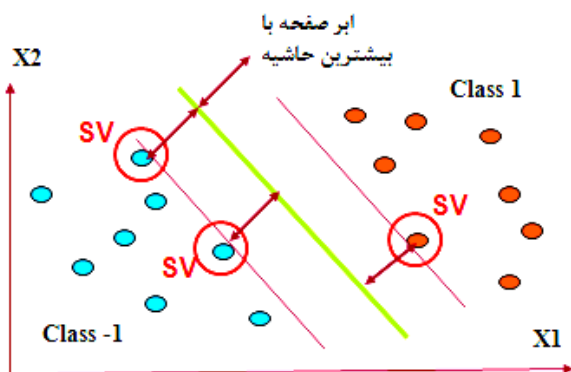
جدول (۱): مثال هایی برای نقش نماهای تقابل استفاده شده.

نقش نمای تقابل	مثال
اما	روایت داستان آنقدر کشش داشت که مخاطب را به سمت خود بکشد و در هر قسمت او را وادار کند پای تلویزیون بنشینند. اما ضعف آن مانند اکثر

مورد نظر ماکزیمم گردد. این نوع انتخاب باعث می شود که تصمیم گیری ما در عمل، شرایط نوبیزی را به خوبی تحمل کند و پاسخ دهی خوبی داشته باشد. این نحوه انتخاب مرز بر اساس نقاطی به نام بردارهای پشتیبان انجام می شود. در SVM با فرض اینکه دسته ها به صورت خطی جداپذیر باشند، ابرصفحه‌هایی با حداکثر حاشیه (Maximum Margin) را بدست می آورد که دسته ها را جدا کنند. برای داده‌های جداپذیر به صورت خطی تابع تفکیک کننده‌گی به شکل زیر نمایش داده می شود:

$$g(x) = \text{sign}(w \cdot \theta(x) + b) \quad (1)$$

در (1) علامت مثبت تعلق به دسته اول و علامت منفی تعلق به دسته دوم را نمایش می دهد، x بردار نمونه‌های آموزشی (ویژگیها)، w بردار ابرصفحه و b باپاس آن می باشد که فاصله بین ابرصفحه تصمیم گیری و مبدا فضای ویژگی را نمایش می دهد. به طور کلی ابرصفحه‌های زیادی برای دسته بندی وجود دارد که در بین آنها SVM، ابرصفحه‌ای را پیدا می کند که فاصله بیشتری را بین دو کلاس ایجاد کند؛ نزدیکترین داده‌های آموزشی به ابرصفحه‌های جداکننده بردار پشتیبان (Support Vector (SV)) نامیده می شوند. این موضوع در شکل 1 نشان داده شده است.



شکل 1: دسته بندی SVM برای جدا کردن یک نوع داده خاص. در این شکل بردارهای پشتیبان و ابرصفحه جداکننده را مشاهده می کنید.

در بسیاری از کاربردهای واقعی ابرصفحه‌ای که بتواند به صورت خطی جداسازی را در دو کلاس انجام دهد وجود ندارد؛ برای حل این مشکل یک راه حل این است که داده‌ها از فضای اصلی به فضای دیگری (فضای با بعد بالاتر) انتقال دهیم و بعد ابرصفحه‌ای را تعیین نماییم که بتواند به صورت خطی این جداسازی را انجام دهد که در این حالت بار محاسباتی زیاد می شود و هزینه محاسبات بالا می رود چرا که به طور کلی ابعاد فضائی که در آن فضا داده‌ها به طور خطی جداپذیر خواهند بود بی نهایت است. از این رو فرآیند نگاشت را توسط توابع هسته

سریالهای ایرانی پایان آن است که به راحتی در چند دیالوگ خلاصه می شود.	
برآورد مبالغ غیر قابل وصول و اندازه گیری درآمد احتمال برابری میان مبالغ مطالباتی که هرساله غیر قابل وصول می گردد و مبلغی که به عنوان ذخیره برای همین رویداد منظور می شود فوق العاده ضعیف است. ولی تفاوت‌های میان این دو رقم اگر به نسبت اندک باشد، قابل چشم پوشی است.	ولی
خانواده مادری یزید، از اعراب مسیحی بودند که بعدا اسلام آوردند. لکن اسلام آنان هیچ گاه آداب و ارزشها و فرهنگ مسیحی آنان را دگرگون نکرده بود.	لکن
یک بار دیگر به بهانه مذاکره با کاسترو درباره خلیج خوکها لباس آغشته به سم را به او هدیه کردند، لیکن کاسترو موضوع را با دیده شک نگریست.	لیکن
در کانادا و آمریکا تنبیه بدنی کودکان توسط والدین آنها تنها در شرایط خاصی جرم محسوب نمی شود. باین حال چنین استثنائی نیز به تدریج در حال رنگ باختن است و گروههای حمایت از کودکان در صدد هستند این موارد را نیز غیرقانونی اعلام کنند.	باین حال
بر روی هر سکه یک مهر خاص ضرب شده بود که نشان ضمانت بهای سکه بود. با وجود این، به نوشته روزنامه ایندپندنت، پس از آن که بر اداره کنندگان کشور روشن شد که از این نوع سکه ها در نقاط دوردست تر پادشاهی استقبال نمی شود، صنعت گران لیدیه به ساخت نخستین کوره ذوب و استخراج طلا همت گماشتند.	باوجود این
سم علف کش اترازین به علت قدرت بسیار زیاد در نابودی علف های هرز و همچنین قیمت پائین آن بسیار پرمصرف است. با این وجود پژوهشگران می گویند باید قوانینی برای محدود ساختن مصرف آن ایجاد کنند.	باین وجود

۵- ماشین بردار پشتیبان

ماشین بردار پشتیبان (Support Vector Machine (SVM)) یکی از انواع روش‌های دسته بندی می باشد که در هر لحظه توانائی دسته بندی دو کلاس را دارد، یا به عبارت دیگر یک نوع دسته بندی دودویی می باشد. در SVM هدف پیدا کردن ابرصفحه‌ای برای جدا کردن داده‌ها در دو کلاس جداگانه است. رویکرد SVM به این صورت است که در فاز آموزش، سعی می شود که مرز تصمیم گیری (Decision Boundary) به گونه‌ای انتخاب گردد که حداقل فاصله آن با هر یک از دسته‌های

به منظور نشان دادن کارایی ویژگی های گفته شده در شناسایی رابطه تقابل به کمک دسته‌بند svm آزمایشاتی صورت گرفت. در ابتدا ۵۰۰۰ نمونه رابطه تقابل که صراحتاً با نقش نماهای خاص رابطه تقابل مانند اما، ولی، لکن، لیکن و ... وجود داشتند، استخراج شد و ۴۰۰۰ نمونه از سایر روابط به تعدادی که در جدول ۲ مشاهده می شود، استخراج گردید؛ همچنین تعداد ۱۰۰۰ نمونه جفت جملاتی را که ما فرض می کنیم رابطه ای بین آنها نیست را استخراج می کنیم. این جفت جملات همجوار نیستند و هر کدام را به طور تصادفی از متون پیکره استخراج می کنیم و سپس جمله اول را به عنوان $Arg1$ و جمله دوم را به عنوان $Arg2$ در نظر می‌گیریم و مانند یک رابطه به استخراج ویژگی می پردازیم.

جدول (۲): تعداد و نوع نقش نماهای استفاده شده.

نقش‌نما	نوع رابطه	تعداد
اما	تقابل	۲۹۰۰
ولی	تقابل	۱۵۰۰
لکن	تقابل	۵۰
لیکن	تقابل	۲۰۰
بالین‌حال	تقابل	۲۵۰
باوجوداین	تقابل	۵۰
بالین وجود	تقابل	۵۰
زیرا	علی	۶۰۰
یعنی	افزایشی	۱۰۰۰
پس	زمانی	۴۰۰
و	افزایشی	۲۰۰
همچنین	افزایشی	۵۰۰
دراین میان	زمانی	۵۰
دراین حال	زمانی	۱۰۰
چون	علی	۱۵۰
بنابراین	علی	۲۰۰
به علاوه	افزایشی	۵۰
برای این که	علی	۵۰
بعد	افزایشی [۲]	۴۰۰
آنگاه	زمانی	۱۵۰
علاوه بر این	افزایشی	۱۵۰

تمامی این نمونه ها از پیکره پردازش هوشمند علائم استخراج شده اند. در این پیکره کلمات یک جمله با توجه به جایگاه نحوی آن کلمه به چند گروه تقسیم شده و برای هر گروه یکسری ویژگی هایی مختص آن گروه بیان شده است برای مثال در این متون برای افعال ویژگی هایی مانند زمان فعل، مرکب یا ساده بودن فعل، شخص فعل، مثبت یا منفی بودن فعل و ... ذکر شده است و یا برای اسامی ویژگی

(Kernel Function) به طور ضمنی انجام خواهیم داد. SVM در داده‌هایی که بصورت خطی جداپذیر نمی‌باشند وابسته به هسته (Kernel) می‌باشد. هسته تابعی است که الگوها را از یک فضای غیرخطی به یک فضای خطی نگاشت می‌کند تا SVM بتواند با استفاده از یک ابرصفحه الگوها را دسته‌بندی کند. در صورت استفاده از توابع هسته، تابع تصمیم‌گیری یک SVM بر پایه هسته به شکل زیر تبدیل می‌گردد:

$$g(x) = \text{sign}\left(\sum_{i=1}^{N_s} \alpha_i y_i k(x_i, x) + b\right) \quad (2)$$

در (۲) $x_i \in R^d, i=1, \dots, N$ ابعاد ویژگیها، N تعداد نمونه‌های آموزشی مربوط به دو دسته، $\alpha_i, y_i \in \{1, -1\}$ ضرایبی که از حل مساله بهینه‌سازی به دست می‌آیند و N_S تعداد بردارهای پشتیبان می‌باشد که مهمترین داده‌های آموزشی هستند چرا که آنها به تنهایی برای تعریف حاشیه تفکیک پذیری کفایت می‌کنند. از معروفترین و متداولترین هسته‌هایی که در SVM به کار برده می‌شوند می‌توان به ۳ و ۴ اشاره کرد:

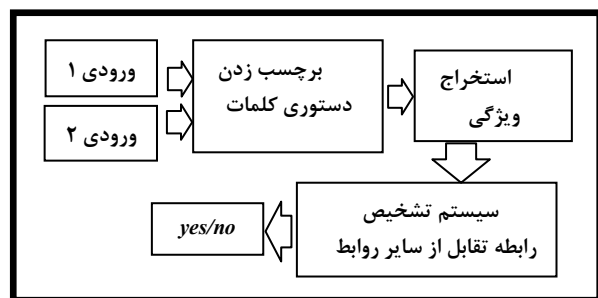
$$k(X, X') = \exp\left(-\frac{1}{2\sigma^2} \|X - X'\|^2\right) \quad (3)$$

$$k(X, X') = (X^T X' + 1)^p \quad (4)$$

که (۳) هسته RBF و (۴) هسته چند جمله‌ای می‌باشد.

۶- سیستم شناسایی رابطه تقابل در گفتمان‌های زبان فارسی

سیستمی که ما قصد طراحی آن را داریم در شکل ۲ نشان داده شده است. در شکل ۲ منظور ما از ورودی ۱ و ورودی ۲ جمله یا بند با پاراگراف قبل و بعد نقش نما با دو عبارت متنی می باشد که پس از برچسب زدن دستوری کلمات که در خود پیکره انجام شده است، به استخراج ویژگی اقدام می‌کنیم و اگر دو ورودی نسبت به هم رابطه تقابل داشته باشند، خروجی yes و اگر نه خروجی no تولید می شود.



شکل ۲- شکل کلی سیستم شناسایی رابطه تقابل

گفتمان می تواند انجام شود و ساخت یک پارسر برای روابط موجود در یک گفتمان و

هایی مانند خاص یا عام بودن اسم، جمع یا مفرد بودن، اسم مکان ذکر گردیده است.

نمونه های استخراج شده را به صورت $\{ [ARG_1] \}$ نقش‌نمای گفتمان $[ARG_2]$ مرتب می کنیم. $([ARG_1] \text{ و } [ARG_2])$ دو قسمت گفتمانی در دوطرف نقش‌نما هستند که می‌تواند بند، جمله یا پاراگراف باشد). در مورد ۱۰۰۰ نمونه بدون رابطه $[ARG_1]$ را جمله اول و $[ARG_2]$ را جمله دوم در نظر می‌گیریم. در ادامه نقش‌نماها را حذف کرده و صرفاً از دو قسمت $[ARG_1]$ و $[ARG_2]$ به عنوان منابعی برای تصمیم‌گیری استفاده می‌کنیم. سپس بردار ویژگی را برای $[ARG_1]$ و $[ARG_2]$ های ایجاد شده برای هر نمونه تشکیل می‌دهیم و نمونه‌ها را به دو کلاس تقسیم می‌کنیم که تمام نمونه های رابطه تقابل را یک کلاس در نظر گرفته و بقیه را کلاس دیگری در نظر می‌گیریم.

در ادامه تمام بردارهای ویژگی که برچسب کلاس مربوطه را دارند را به نرم افزار $MATLAB 7.8$ منتقل کرده داده‌ها را به داده‌های آموزش و آزمایش تقسیم کرده و به منظور آموزش دسته‌بند SVM استفاده می‌کنیم که مراحل آموزش و آزمایش را با توابع هسته مختلف امتحان می‌کنیم که نتایج و میزان صحت را در جدول ۳ مشاهده می‌کنید.

جدول ۳: درصد صحت روش پیشنهادی با دسته‌بند SVM و توابع هسته

مختلف

نوع هسته	درصد صحت
خطی	۶۱.۰۲
درجه ۲	۶۰.۵۹
درجه ۳	۵۹.۷۴
RBF	۶۰.۹۶

۸- مراجع

- [۱] ذوقدارمقدم، رضا و دبیرمقدم، محمد، "نقش‌نماهای گفتمان مقایسه نقش‌نمای *but* در زبان انگلیسی با نقش‌نمای اما در زبان فارسی"، پژوهش زبان‌های خارجی، ۱۲، ۵۵-۷۶، ۱۳۸۱.
- [۲] مقدم‌کیا، رضا، "بعد، نقش‌نمای گفتمان در زبان فارسی"، نامه فرهنگستان، ۲۳، ۸۱-۹۸، ۱۳۸۴.
- [۳] D. Marcu and A. Echiabi. ۲۰۰۱. *An unsupervised approach to recognizing discourse relations. In Proceedings of the ۴۰th Annual Meeting on Association for Computational Linguistics, pages ۳۶۸-۳۷۵.*
- [۴] Manami Saito, Kazuhide Yamamoto, and Satoshi Sekine. ۲۰۰۶. *Using phrasal patterns to identify discourse relations. In Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics (HLT-NAACL ۲۰۰۶), pages ۱۳۳-۱۳۶, New York, USA, June.*
- [۵] Emily Pitler, Annie Louis, and Ani Nenkova. ۲۰۰۹. *Automatic sense prediction for implicit discourse relations in text. To appear in Proceedings of the Joint Conference of the ۴۷th Annual Meeting of the Association for Computational Linguistics and the ۴th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing (ACL-IJCNLP ۲۰۰۹), pages ۶۸۳-۶۹۱.*
- [۶] Z.H. Lin, M.Y. Kan and H.T. Ng. ۲۰۰۹. *Recognizing Implicit Discourse Relations in the Penn Discourse Treebank. Proceedings of the ۲۰۰۹ Conference on Empirical Methods in Natural Language Processing (ACL and AFNLP ۲۰۰۹), pages ۳۴۳-۳۵۱.*
- [۷] Zhou Z, M Lan, X Yu, Z Niu, J Su and C L Tan, *Predicting discourse connectives for implicit discourse relation recognition, ۲۳rd International Conference on Computational Linguistics (COLING ۲۰۱۰), August ۲۳-۲۷, ۲۰۱۰, Beijing, China.*
- [۸] Rashmi Prasad, Nikhil Dinesh, Alan Lee, Eleni Miltsakaki, Livio Robaldo, Aravind Joshi, and Bonnie Webber. ۲۰۰۸. *The Penn Discourse Treebank ۲.۰. In Proceedings of the ۶th International Conference on Language Resources and Evaluation (LREC ۲۰۰۸).*
- [۹] E. Pitler and A. Nenkova. ۲۰۰۹. *Using Syntax to Disambiguate Explicit Discourse Connectives in Text. Proceedings of the ACL-IJCNLP ۲۰۰۹ Conference Short Papers.*

۷- نتیجه‌گیری و کارهای آینده

در این مقاله روشی برای شناسایی رابطه تقابل در گفتمان‌های فارسی ارائه شد که اولین کار در حوزه پردازش کامپیوتری گفتمان می‌باشد. از آنجا که در رابطه تقابل نیز ممکن است در مواردی نقش‌نما حذف شود، بنابراین سیستمی مستقل از نقش‌نما ارائه شد که با استخراج ویژگی‌های سعی در شناسایی رابطه تقابل دارد. نتایج دسته‌بند SVM بر روی ویژگی‌های استخراج شده، امیدوارکننده بود و در بهترین حالت بیشتر از ۶۰ درصد صحت ایجاد شد. به عنوان پیشنهاد کار برای آینده می‌توان ویژگی‌های بهتری را استخراج کرد مانند پیدا کردن کلمات متضاد در دو طرف رابطه و یا پیدا کردن الگوها و عباراتی که نقش‌نما نیستند اما در رابطه تقابل ممکن است به کار گرفته شوند. به علت این که تحقیقات زیادی در پردازش‌های کامپیوتری گفتمان در زبان فارسی انجام نشده است، کارهای زیادی را می‌توان در این حوزه انجام داد. به عنوان مثال شناسایی سایر روابط موجود در یک