

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

١١٤٩ ٣٣



دانشگاه اصفهان

دانشکده فنی و مهندسی

گروه مهندسی کامپیووتر

پایان نامه کارشناسی ارشد رشته مهندسی کامپیووتر گرایش نرم افزار

رفع ابهام معنای کلمات در حاشیه نویسی خودکار اسناد توسط ادغام روش یادگیری ماشین و روش مبتنی بر دانش

استاد راهنما:

دکتر کامران زمانی فر

پژوهشگر:

بهنام حاجیان

اعلامات مذکون چشمی برداشته
مذکون

اردیبهشت ماه ۱۳۸۸

۱۳۸۸/۴/۶

۱۱۴۹۲۳

کلیه حقوق مادی مرتب بر نتایج مطالعات، ابتكارات
و نوآوری های ناشی از تحقیق موضوع این پایان نامه
متعلق به دانشگاه اصفهان است.

پیووه کارشناسی پایان نامه
دراحت شده است
تخصیلات تکمیلی دانشگاه اصفهان



دانشگاه اصفهان

دانشکده مهندسی

گروه مهندسی کامپیوتر

پایان نامه کارشناسی ارشد رشته مهندسی کامپیوتر گرایش نرم افزار
آقای بهنام حاجیان تحت عنوان

رفع ابهام معنای کلمات در حاشیه نویسی خودکار اسناد توسط ادغام روش
یادگیری ماشین و روش مبتنی بر دانش

در تاریخ ۱۳۹۰/۰۷/۲۸ توسط هیئت داوران زیر بررسی و با درجه ممتاز... به تصویب نهایی رسید.

۱- استاد راهنمای پایان نامه: دکتر کامران زمانی فر با مرتبه علمی استادیاری امضا.....

۲- استاد داور داخل گروه: دکتراحمد برا آنی با مرتبه علمی استادیاری امضا.....

۳- استاد داور خارج از گروه: دکتر مازیار پالهنج با مرتبه علمی استادیاری امضا.....

امضاء مدیر گروه
امیرنژاد
۱۴۰۰/۰۷/۲۹

بسم الله الرحمن الرحيم

سپاسگزاری

سپاس خدای یکتا را که هر چه هست از اوست.

سپاس او را که بی نام و یادش نمی توان آغاز کرد و نمی توان به پایان برد.

سپاس و حمد بیکران او را که به ما طریق عاشقی آموخت.

سپاس ترا ای یگانه عالم، ای حقیقت مطلق، ای منتهای تحقیق، ای همه هستی در یک کلام، ای خدا

خدایا مخواه که غرور در دل ما ریشه کند

خدایا مخواه که غیر از آنچه تو خواسته ای بخواهیم و غیر از آنچه رضای توست بطلبیم

خدایا ما را آنی به خود و امگذار که به خود نیامده ایم که به خود رویم

خدایا تو میدانی که جز تو از تو نمی خواهیم

تنها تو را می پرستیم و تنها از تو یاری می طلبیم

ما را در این راه تنها مگذار

به مصدق کلام شریف من لم یشکر المخلوق لم یشکر الخالق، بر خود فرض میدانم تا از تلاش تمامی راهنمایان خود در این پروژه بویژه جناب آقای دکتر زمانی فر تشرکتم و از پروردگار یکتا توفیق روز افزون ایشان را در خدمت به ایران اسلامی مسئلت دارم.

از سروران و عزیزانی که با رهنمودهای خود ما را در یافتن کاستی ها و رفع آنها یاری خواهند داد سپاسگزارم و امیدوارم که شایستگی پذیرش نظرات و پیشنهادات این عزیزان را داشته باشم.

من الله التوفيق

بهنام حاجیان

تقدیم به:

پدر و مادرم که در تمام مراحل زندگی مرا
راهنمایی و پشتیبانی کرده اند.

چکیده

با ظهور وب معنایی به عرصه اینترنت، محیطی قابل درک و استنتاج برای محاسبات کامپیوتری فراهم گردید. در این زمینه اولین گام جهت دستیابی به محیط هوشمند و قابل فهم برای ماشین، ایجاد هستی شناسی مناسب در دامنه مورد بحث و اتصال اطلاعات موجود در صفحات وب به نمونه ها و کلاس های موجود در هستی شناسی می باشد. این عمل که حاشیه نویسی معنایی اسناد نام دارد، کلیه کلمات ، اشیاء و اطلاعات موجود در صفحات وب را به مفاهیم متضایر آنها در هستی شناسی مرتبط می سازد لذا دامنه معنایی هر کلمه با توجه به کلاس مربوطه مشخص می گردد. بدین جهت بسیاری از مشکلات مطرح شده در زمینه بازیابی و جستجوی اطلاعات حل شده و اطلاعات موجود در اسناد به مفاهیمی قابل فهم برای کامپیوترها تبدیل می شوند. در این زمینه خصوصاً با توجه به وجود کلمات هم شکل و با معانی متفاوت، با محیطی مبهم روبرو هستیم. لذا معنای کلمات می بایست با توجه به زمینه و مفهوم جمله ای که در آن بکار رفته است تعیین و رفع ابهام گردد. عمل حاشیه نویسی معنایی اسناد در حال حاضر به طور دستی و یا نیمه خودکار توسط عامل انسانی انجام می گیرد که عملی وقتگیر، پیچیده، پرهزینه و نادقیق می باشد.

یکی از مسائل زیر مجموعه حاشیه نویسی خودکار اسناد رفع ابهام معنی کلمات در متن می باشد. در این پروژه سعی شده با ارائه راهکاری بر مبنای هوش مашینی و بکارگیری الگوریتم های یادگیری ماشینی با ناظر به خودکار نمودن این عمل وقتگیر توسط کامپیوتر پرداخته شود. سیستم فوق دارای دو فاز آموزش و عملیاتی می باشد. در مرحله آموزش سیستم توسط داده های آموزشی (صفحات حاشیه نویسی شده بصورت دستی) آموزش دیده و سپس در فاز عملیاتی به حاشیه نویسی اسناد جدید بصورت خودکار می پردازد. در این پایان نامه از چند الگوریتم مختلف آموزش ماشین استفاده شده است و نتایج عملکرد آنها ارزیابی و با یکدیگر مقایسه شده است. در این ارتباط الگوریتم ماشین بردار پشتیبان بهترین کارایی را از خود نشان داده است. در مرحله بعد جهت افزایش دقت و کارایی این الگوریتم ها به ارائه چند مرحله پیش پردازش پرداخته شده است. نتایج ارزیابی در این پایان نامه نشان خواهد داد که با ادغام روش های بر مبنای یادگیری ماشینی و سیستم های مبتنی بر دانش می توان به بهترین کارایی دست پیدا نمود. همچنین استفاده از روش انتخاب ویژگی ها در این مسئله از دقت سیستم می کاهد. در صورتی که بوسیله گسترش ویژگی های متن، دقت و کارایی سیستم افزایش خواهد یافت.

کلید واژه ها:

رفع ابهام معنایی ، وب معنایی ، حاشیه نویسی وب ، یادگیری ماشینی ، هستی شناسی ، متن کاوی ، گسترش ویژگی ها

فهرست مطالب

صفحه	عنوان
۱	فصل ۱ : معرفی
۲	۱-۱ مقدمه
۳	۲-۱ تعریف مسئله و اهداف پایان نامه
۴	۱-۲-۱ هدف پایان نامه
۴	۳-۱ روش انجام پژوهش
۵	۴-۱ ساختار پایان نامه
۷	فصل ۲ : مفاهیم وب معنایی و یادگیری ماشینی
۷	۱-۲ مقدمه
۹	۲-۲ تعریف وب معنایی
۱۳	۳-۲ پیش نیازهای وب معنایی
۱۳	۱-۳-۲ هستی شناسی
۱۵	۱-۱-۳-۲ انواع هستی شناسی
۱۶	۲-۱-۳-۲ کاربردهای هستی شناسی
۱۷	۳-۱-۳-۲ زبانهای هستی شناسی
۱۸	۱-۳-۱-۳-۲ XML
۱۹	۲-۳-۱-۳-۲ RDF : چهارچوب توصیفی منابع
۲۰	۳-۱-۳-۲ OWL : زبان هستی شناسی وب
۲۱	۴-۱-۳-۲ نحوه ایجاد هستی شناسی و مهندسی هستی شناسی
۲۲	۲-۳-۲ فرا داده در وب و حاشیه نویسی وب معنایی
۲۴	۳-۳-۲ جایگاه منطق در وب معنایی
۲۵	۴-۲ عاملها
۲۶	۵-۲ مسائل قابل حل توسط وب معنایی
۲۶	۱-۵-۲ مدیریت دانش
۲۸	۲-۵-۲ تجارت الکترونیک
۲۹	۳-۵-۲ سرویس‌های جستجو

صفحه	عنوان
۳۲	۶-۲ مشکلات مطرح موجود در وب معنایی.....
۳۳	۷-۲ مفاهیم داده کاوی و یادگیری ماشینی.....
۳۳	۱-۷-۲ مقدمه ای بر داده کاوی.....
۳۵	۲-۷-۲ مدل‌های داده کاوی.....
۳۵	۱-۲-۷-۲ رده بندی
۳۵	۲-۲-۷-۲ خوشه بندی
۳۶	۳-۲-۷-۲ تحلیل پیوند (تحلیل وابستگی داده ها).....
۳۶	۳-۷-۲ تعریف یادگیری ماشینی
۳۷	۴-۷-۲ معرفی برخی از الگوریتمهای یادگیری ماشین و داده کاوی
۳۷	۱-۴-۷-۲ شبکه های عصبی
۳۸	۲-۴-۷-۲ درختهای تصمیم
۳۸	۳-۴-۷-۲ استنتاج قوانین
۳۸	K4-۴-۷-۲ نزدیکترین همسایه
۳۹	۵-۴-۷-۲ الگوریتم کلاس بند ناییوبیز
۴۰	۶-۴-۷-۲ الگوریتم ماشین بردار پشتیبان
۴۱	۵-۷-۲ انتخاب ویژگی ها در فرایند یادگیری ماشینی
۴۲	۸-۲ جمع بندی
فصل ۳ : کارهای انجام شده پیشین در زمینه حاشیه نویسی معنایی و رفع ابهام معنی کلمات	
۴۳	۱-۳ مقدمه
۴۴	۲-۳ تقسیم بندی کارهای انجام شده در زمینه حاشیه نویسی معنایی وب و یادگیری هستی شناسی ...
۴۶	۳-۳ دسته بندی سکوها و ابزارهای حاشیه نویسی
۴۶	۱-۳-۳ ابزارهای حاشیه نویسی دستی
۴۸	۲-۳-۳ ابزارهای حاشیه نویسی خودکار و نیمه خودکار
۴۸	۱-۲-۳-۳ روشهای حاشیه نویسی مبتنی بر الگو
۴۹	۲-۲-۳-۳ روشهای حاشیه نویسی خودکار مبتنی بر یادگیری ماشینی
۴۹	۴-۳ نگاهی به سکوها و سیستم های حاشیه نویسی معنایی
۵۰	۱-۴-۳ سیستم Aero DAML
۵۱	۲-۴-۳ سکوی آرمادیلو

صفحة	عنوان
۵۱	۳-۴-۳ سکوی حاشیه نویسی مدیریت دانش و اطلاعات (KIM)
۵۲	۴-۴-۳ سیستم آن ام
۵۳	۴-۴-۳ ابزار موسه
۵۴	۴-۴-۳ سیستم آنتومت
۵۴	۴-۴-۳ سیستم سمتگ و سکوی سیکر
۵۵	۴-۴-۳ سیستم آمیلکیر
۵۶	۴-۴-۳ سیستم اسکریم
۵۶	۴-۴-۳ سیستم مجی پای
۵۶	۴-۴-۳ سیستم حاشیه نویسی سی پنکو
۵۹	۴-۳-۳ دسته بندی کارهای انجام شده در زمینه رفع ابهام معنای کلمات
۶۰	۴-۳-۳ جمع بندی
۶۱	فصل ۴ : طراحی سیستم پیشنهادی جهت رفع ابهام معنای کلمات موجود در متن
۶۱	۴-۱ مقدمه
۶۴	۴-۲ تعریف مسئله
۶۵	۴-۳-۴ پیش نیاز های ساخت سیستم حاشیه نویسی خودکار وب
۶۵	۴-۴-۴ استاندارد های حاشیه نویسی و زبان مناسب بیان فرا داده ها
۶۶	۴-۳-۴ هستی شناسی مرتبط با دامنه مورد بحث
۶۹	۴-۳-۴ داده های آموزشی و ابزار ذخیره سازی و بازیابی اطلاعات
۶۹	۴-۴-۴ معماری سیستم حاشیه نویسی معنایی خودکار سوام
۷۰	۴-۴-۴ شرح معماری سیستم سوام
۷۳	۴-۴-۴ مراحل پیش پردازش
۷۳	۴-۲-۴-۴ ۱- جداسازی و شماره گذاری کلمات و جملات موجود در متن
۷۳	۴-۲-۴-۴ ۲- حذف کلمات اضافه و اعداد از متن
۷۴	۴-۲-۴-۴ ۳- پیش پردازش کلمات متراff با کمک پایگاه دانش ووردن (گسترش ویژگی ها)
۷۴	۴-۲-۴-۴ ۴- ریشه یابی کلمات موجود در متن
۷۴	۴-۲-۴-۴ ۵- پیش پردازش و ساخت بردار فضای حالت سند
۷۵	۴-۲-۴-۴ ۶- اعمال صافی عکس بسامد تکرار کلمه در متن ، نرمال سازی و انتخاب ویژگی ها
۷۶	۴-۴-۴ مرحله آموزش

عنوان	
صفحه.	
۷۶	۴-۴ مرحله حاشیه نویسی معنایی اسناد.....
۷۸	۵-۴ جمع بندی فصل.....
۷۹	فصل ۵ : پیاده سازی و ارزیابی نتایج سیستم سوام
۷۹	۱-۵ مقدمه.....
۸۰	۲-۵ پیاده سازی سیستم سوام
۸۰	۱-۲-۵ مولفه ذخیره و بازیابی اطلاعات و اسناد
۸۱	۲-۲-۵ مولفه پیش پردازش
۸۵	۳-۲-۵ مولفه یادگیری ماشینی الگوریتم های کلاس بندی متداول.....
۸۵	۴-۲-۵ مولفه الگوریتم کلاس بندی ACC
۹۰	۵-۲-۵ مولفه حاشیه نویسی
۹۱	۶-۲-۵ مولفه کار با هستی شناسی شبکه واژگان
۹۱	۷-۲-۵ مولفه ارزیابی سیستم.....
۹۲	۳-۵ معرفی بستر مورد ارزیابی سیستم سوام
۹۶	۴-۵ نتایج حاصل از ارزیابی سیستم
۱۰۳	۵-۵ بررسی نتایج حاصل از آزمایشات ارزیابی سیستم.....
۱۱۳	۶-۵ جمع بندی.....
۱۱۴	فصل ۶ : تحلیل نتایج ، نتیجه گیری و کار های آینده
۱۱۴	۱-۶ مقدمه.....
۱۱۵	۲-۶ تحلیل نتایج بدست آمده در پایان نامه
۱۱۷	۳-۶ نتیجه گیری
۱۱۹	۴-۶ پیشنهادات و کار های آینده
۱۲۱	پیوست ۱ : شمای واسط کاربر ارزیابی نتایج سیستم
۱۲۲	واژه نامه (انگلیسی به فارسی)
۱۲۴	منابع و مأخذ

فهرست شکل ها

عنوان	صفحة
شکل ۱-۲ سطوح مختلف قابلیت ادراک در داده	۱۰
شکل ۲-۲ طرح لایه ای وب معنایی [۱۳]	۱۲
شکل ۳-۲ نمونه ساختاریک فایل XML	۱۸
شکل ۴-۲ نمونه نحو زبان RDF برای نمایش سه تایی Buddy Belden Owns a Business	۱۹
شکل ۵-۲ نمایش RDF جمله "the people at the meeting were joe, Bob, Susan and Ralph"	۲۰
شکل ۶-۲ نمایش بصری هستی شناسی ذکر شده در شکل ۵-۲	۲۰
شکل ۷-۲ مراحل ایجاد هستی شناسی و مهندسی هستی شناسی [۱۳]	۲۱
شکل ۸-۲ شمای حاشیه نویسی معنایی [۱۸]	۲۳
شکل ۹-۲ نمونه ای از فرا داده های معنایی به زبان XML	۲۳
شکل ۱۰-۲ نحوه عملکرد الگوریتم ماشین بردار پشتیبان.	۴۱
شکل ۱-۳ فرایند حاشیه نویسی دستی استاد وب	۴۷
شکل ۲-۳ فرایند حاشیه نویسی خودکار	۴۸
شکل ۳-۳ فرایند انجام شده در سیستم پنکو [۱۹]	۵۷
شکل ۱-۴ معماری سیستم حاشیه نویسی سوام SWAM	۷۱
شکل ۲-۵ مرحله استخراج مجموعه های کاندید تکرار شونده در ACC [۴۸]، [۳۲]	۸۸
شکل ۳-۵ مرحله استخراج قوانین پیوند در الگوریتم ACC [۴۸]	۸۸
شکل ۴-۵ نمونه ای از قوانین استخراج شده توسط الگوریتم ACC در مرحله آموزش.	۸۹
شکل ۵-۵ شمای سیستم حاشیه نویسی سوام به همراه ابزار گیت.	۹۰
شکل ۶-۵ فرمت ذخیره یک نمونه از اخبار مجموعه داده رویتر	۹۲
شکل ۷-۵ نمونه مجموعه داده های موجود در کورپوس سنس ایوال	۹۵
شکل ۸-۵ تاثیر پارامتر حداقل درجه حمایت بر معیار های فراخوانی و دقت الگوریتم ACC	۱۰۴
شکل ۹-۵ تاثیر پارامتر حداقل درجه اطمینان بر معیار های فراخوانی و دقت الگوریتم ACC	۱۰۵
شکل ۱۰-۵ تاثیر اضافه نمودن واژه های فرا معنی بر دقت سیستم سوام.	۱۰۷
شکل ۱۱-۵ توزیع دقت عملکرد سیستم در کلمات مختلف موجود در دیتا است سنس ایوال	۱۱۰
شکل ۱۲-۵ رابطه میان دقت سیستم و تعداد نمونه های موجود در مجموعه داده مورد آزمایش	۱۱۱
شکل ۱۳-۵ روند اثر تعداد معانی معرفی شده در مجموعه داده های سنس ایوال در دقت سیستم	۱۱۲

فهرست جدول ها

عنوان	صفحه
جدول ۱-۳ جمع بندی خصوصیات ابزارهای حاشیه نویسی [۲۰].	۵۸
جدول ۱-۴ تعداد کلمات موجود در پایگاه دانش هستی شناسی شبکه واژگان در سال ۲۰۰۶ [۴۲].	۶۷
جدول ۲-۴ انواع روابط مربوط به اسم های موجود در سیستم شبکه واژگان (ووردنت) [۹].	۶۸
جدول ۳-۴ انواع روابط مربوط به فعل های موجود در سیستم شبکه واژگان (ووردنت) [۹].	۶۹
جدول ۱-۵ لیست کلمات توقف و کلمات اضافه که از متن حذف می گردند.	۸۲
جدول ۲-۵ توزیع کلمات موجود در مجموعه داده سنss ایوال از نظر نقش کلمه [۴۶].	۹۴
جدول ۳-۵ تاثیر پارامتر حداقل درجه حمایت بر پوشش و دقت الگوریتم ACC	۹۶
جدول ۴-۵ تاثیر پارامتر حداقل درجه اطمینان بر پوشش و دقت الگوریتم ACC	۹۷
جدول ۵-۵ نام و گروه طبقه بندی الگوریتم های استفاده شده در سیستم سوام.	۹۸
جدول ۶-۵ نتیجه ارزیابی دقت سیستم بدون اعمال مراحل پیش پردازش.	۹۸
جدول ۷-۵ تاثیر اضافه نمودن مرحله افزودن واژه های مترادف هایپرنیم بر دقت سیستم.	۹۹
جدول ۸-۵ تاثیر اضافه نمودن مرحله حذف واژه های توقف بر کارایی سیستم.	۹۹
جدول ۹-۵ تاثیر اضافه نمودن مرحله ریشه یابی بر کارایی سیستم.	۱۰۰
جدول ۱۰-۵ تاثیر اضافه نمودن صافی TFIDF بر کارایی سیستم.	۱۰۰
جدول ۱۱-۵ تاثیر اضافه نمودن انتخاب ویژگی بوسیله صافی بهره اطلاعاتی بر دقت سیستم.	۱۰۰
جدول ۱۲-۵ نتایج ارزیابی دقت سیستم توسط الگوریتم ماشین بردار پشتیبان بر داده های سنss ایوال.	۱۰۱
جدول ۱۳-۵ اثر اعمال مراحل ترکیبی پیش پردازش بر روی کارایی الگوریتم ماشین بردار پشتیبان در رفع ابهام معنایی کلمات در نقش اسم.	۱۰۳
جدول ۱۴-۵ مقایسه بهترین الگوریتم ها پس از اعمال مراحل (اضافه نمودن هایپرنیم + ریشه یابی + حذف واژگان توقف) بر روی کلمات در نقش اسم.	۱۰۳
جدول ۱۵-۵ میانگین دقت سیستم توسط الگوریتم ماشین بردار پشتیبان با اعمال مراحل پیش پردازش، بر روی کلمات در نقش های مختلف.	۱۰۸
جدول ۱۶-۵ نتایج رفع ابهام تیم های شرکت کننده در مسابقه سنss ایوال ۲ در رفع ابهام کلمات در نقش اسم.	۱۰۸

فهرست کوتاه نوشت ها

XML	Extensible Markup Language
RDF	Resource Description Framework
OWL	Ontology Web Langiage
SVM	Support Vector Machine
SWAM	Semantic Wen Annotation by Machine Learning
WSD	Word Sense Disambiguation
ACC	Association Concept Classifier
UML	Unified Markup Language
HTML	Hyper Text Markup Language
FOL	First Order Logic
W³C	Word Wide Web Consortium
DAML	Darpa Agent Markup Language
B²B	Business to Business
B²C	Business to Commerce
API	Application Programming Interface
IDE	Integrated Development Environment
POS	Ports of Speech
VSM	Vector Space Model

خ

فصل ۱

معرفی

امروزه با توجه به گسترش روز افزون هوش مصنوعی در علم کامپیوتر، بسیاری از عملیات دستی جای خود را به نرم افزار های هوشمند کامپیوتری داده است. بطور کلی ماهیت وجودی هوش به مفهوم جمع آوری اطلاعات، استقراء و تحلیل تجربیات به منظور رسیدن به دانش و یا ارائه تصمیم می باشد. هوش مصنوعی، علم و مهندسی ماشین های هوشمند با بکارگیری از کامپیوتر و الگوگیری از درک هوش انسانی و نهایتاً دستیابی به مکانیزم هوشمند می باشد [۱]. در سالهای اخیر این هوشمندی کاملاً در عامل های نرم افزاری خود مختار و هوشمند نمایان شده است بطوریکه این عامل ها در بسیاری از موارد جای انسان ها را گرفته و وظایف آنها را انجام می دهند. به عنوان مثال این ابزار می تواند به عنوان یک کارپرداز از کاربران دستور گرفته و پس از جستجو در سایت ها اقدام به و رزرو و خرید بلیط هوا پیما با قیمت مناسب بکند و یا در فعالیت های تجاری مثل یک عامل انسانی اقدام به خرید و فروش و یا معامله نماید. جهت دستیابی به چنین دنیای هوشمندی، نیاز به زیر ساخت هایی جهت تبادل اطلاعات، ابزار های جستجو و دانش زمینه جهت استنتاج می باشد.

۱-۱ مقدمه

امروزه وب به عنوان یکی از مهمترین بسترهای جمع آوری اطلاعات برای عامل‌های هوشمند معرفی می‌شود. لذا یکی از دغدغه‌های دنیای کامپیوتر، ایجاد بستری مناسب در وب جهت ایجاد امکان تعامل بین عامل‌های هوشمند و دنیای وب می‌باشد. متاسفانه اطلاعات موجود در وب متداول، قابل پردازش و درک برای این نرم افزار‌های هوشمند نیست. این اطلاعات کاملاً بصورت غیر ساخت یافته در دنیای اینترنت گسترش یافته‌اند. اغلب این اطلاعات بصورت متنی بوده و به زبان طبیعی، در دنیای وب در دسترس کاربران قرار گرفته‌اند. ولی زبان‌های طبیعی غیر قابل فهم و پردازش برای نرم افزار‌های کامپیوتری می‌باشند. لذا عامل‌های هوشمند توانایی تعامل با وب کنونی را نخواهند داشت [۳]، [۲]. از سوی دیگر، با توجه به حجم بسیار زیاد اطلاعات موجود در وب مشکلات عمده‌ای در زمینه بازیابی و جستجوی اطلاعات مورد نظر کاربران پدیدار شده است. در این حیطه موتور‌های جستجو نیز قادر به ارائه نتایج دقیق و ارزشمندی برای کاربران نمی‌باشند. موتور‌های جستجوی متداول تنها قادر به جستجوی استناد با توجه به وجود یا عدم وجود کلید واژه‌های موجود در پرس و جوی وارد شده توسط کاربر بوده و در کنار آن تنها می‌توانند استناد جدید را که توسط خزشگر^۱ یافت می‌شوند به پایگاه داده شان افزوده و یا صفحات قبلی را به روز رسانی نمایند. وجود حجم وسیع اطلاعات در وب و ازدیاد صفحات وب^۲ و فقدان عقل سلیم و عدم توانایی موتور‌های جستجو از تحلیل معنایی صفحات باعث گردیده که نتایج جستجو در این ابزار برای کاربران چندان دقیق و قابل قبول نباشد و امروز شاهد وجود بسیاری از صفحات نا مرتب با پرس و جوی کاربر در نتایج این ابزار می‌باشیم [۴]، [۲].

با توجه به وجود نا رسانی هایی در وب متداول، پس از ارائه مقاله "تیم برنزیلی" در سال ۲۰۰۱ میلادی، وب معنایی پا به عرصه‌ی دنیای اینترنت گذاشت [۵]. این محقق و مخترع وب، اهداف و دید دنیای وب را فراتر از صورت متداول آن در دنیای امروز می‌داند بطوریکه اطلاعات موجود در وب معنایی را قابل پردازش و تحلیل توسط ماشین می‌خواند [۶]. با وجود آمدن زبان XML قادر به نوعی طبقه‌بندی ساختاری در این اطلاعات شده ولی هنوز این اطلاعات قابل درک و استنتاج توسط کامپیوتر‌ها نبوده و هنوز خلاء معنایی در ساختار اطلاعات موجود در وب احساس می‌شد. وب معنایی این خلا را توسط حاشیه نویسی استناد وب و توسعه وب جاری پر

^۱ Crawler

^۲ ۳۰ میلیارد صفحه در فوریه سال ۲۰۰۷ طبق آمار نت کرافت

کرده که در آن اطلاعات معنایی بصورت خوش تعریف در قسمتی از پایگاه دانش^۳ کامپیوتر که هستی شناسی^۴ نام دارد تعریف میگردد و قابلیت استنتاج از اطلاعات موجود در وب را به عاملهای هوشمند اهدا نموده است. همراه با این پیشرفت، عامل های هوشمند می توانند با بهره گیری از منطق و استنتاج از قوانین و حقایق تعریف شده در پایگاه دانش و اطلاعات موجود در هستی شناسی و اتصال اطلاعات موجود در وب به مفاهیم تعریف شده در آن اقدام به استنتاج و جوابگویی مناسب به درخواست ها و پرسش های کاربر نمایند[۳].

همچنین توسط عمل حاشیه نویسی بسیاری از مشکلات موجود در موتور های جستجو با توجه به تعیین و وضوح معنای دقیق برای کلمات حل می گردد. زیرا در وب معنایی، موتور های جستجو می توانند به معنای کلمات توسط فرا داده هایی که معنای کلمات را با توجه به کلاس هستی شناسی دامنه مورد بحث معرفی می کنند پی بيرند و اقدام به حذف صفحات نا مرتبط که شامل کلید واژه های پرس و جو با معنای غیر از معنای مورد نظر کاربر است نمایند. در نتیجه قادر به نمایش نتایج دقیق تری به کاربران می باشد.

عمل حاشیه نویسی و ساخت هستی شناسی در حال حاضر به صورت دستی و یا ابزار نیمه خودکار انجام می گیرد که عملی وقتگیر و پر هزینه می باشد. این مشکل برای وب سایت های بسیار بزرگ مثل وب سایت های خبری بصورت دستی غیر قابل انجام بوده و نیاز به ابزار هایی جهت خودکار نمودن تولید این فرا داده ها احساس می گردد.

۲-۱ تعریف مسئله و اهداف پایان نامه

با توجه به مشکلات مطرح شده در بالا و پیشنهاد حل آنها توسط وب معنایی، حاشیه نویسی معنایی اسناد وب یکی از مهمترین دغدغه های مطرح شده در ایجاد این ساختار در وب می باشد. حاشیه نویسی معنایی اسناد به معنی اضافه نمودن فرا داده هایی در وب به منظور تعیین معنای اطلاعات موجود در وب و افزودن مفاهیم قابل پردازش و فهم توسط ماشین می باشد [۷]. همانگونه که در بالا ذکر شد این عمل در حال حاضر در محدوده کوچکی از منابع توسط عامل انسانی انجام می گیرد. مسئله تبدیل وب سایتهای جاری به وب سایت های معنایی تنها توسط عمل حاشیه نویسی قابل انجام است که می بایست توسط طراحان و تولید کنندگان وب سایت ها انجام گیرد. مسئله ساخت دستی نمونه ها در هستی شناسی و حاشیه نویسی صفحات نیاز به تخصص در مهندسی هستی

^۳ Knowledge Base

^۴ Ontology

شناسی و چگونگی ساخت و نگهداری آن دارد. از طرف دیگر انجام حاشیه نویسی بصورت دستی دقیق نمی باشد و نیاز به اتلاف زمان و هزینه بسیار دارد [۸]. یکی از مسائل مهم در مسئله حاشیه نویسی معنایی اسناد، ارائه روشی با دقت کافی جهت رفع ابهام معنای کلمات هم شکل با معنی متفاوت^۰ است که هدف اصلی این پایان نامه را تشکیل می دهد.

۱-۲-۱ هدف پایان نامه

با توجه به دقت مورد نیاز در رفع ابهام کلمات^۱ موجود در متن، هدف اصلی در این پایان نامه یافتن روشی کارا و با دقت کافی جهت تشخیص رفع ابهام معنای کلمات در جمله ای است که در آن ظاهر شده اند. به عبارت دیگر در این پایان نامه پس از ارائه معماری مناسبی برای حاشیه نویسی اسناد بر مبنای یادگیری ماشین به ارائه راهکار هایی جهت بهبود الگوریتم های رده بندی بوسیله روش های مبتنی بر دانش جهت حل مسئله ای رفع ابهام معنی کلمات پرداخته شده است.

۳-۱ روش انجام پروژه

در این پایان نامه پس از بررسی کارهای انجام شده در گذشته بر روی مسئله فوق، سیستمی ترکیبی بر اساس مدل یادگیری ماشین^۲ ارائه شده تا بتواند اطلاعات متى موجود در صفحات وب را با توجه به معنای صحیح برآمده از کلمه در هر جمله به مفهوم صحیح و کلاس معنایی مقتضی مربوطه در هستی شناسی اختصاص دهد. به عبارت دیگر سیستم فوق قادر است معنای صحیح کلمات استفاده شده در متن را بوسیله آموزش از صفحاتی که قبل بتصورت دستی حاشیه نویسی شده اند تشخیص داده و به کلاس صحیح موجود در هستی شناسی انتساب داده و حاشیه نویسی کند. سپس با ترکیب روش یادگیری ماشینی و روش مبتنی بر دانش^۳ کارایی و دقت سیستم تا حد قابل قبولی بالا برد شده است.

^۰ Polysemy

^۱ Word Sense Disambiguation

^۲ Machine Learning

^۳ Knowledge Based

گفته است سیستم فوق به ایجاد نمونه ها در هستی شناسی نمی پردازد، بلکه سیستم فوق جهت تطابق نمونه های موجود در هستی شناسی شبکه واژگان^۹ با کلمات موجود در متن ساخته و ارزیابی شده است، به گونه ای که بتواند در معنی کلمات هم شکل رفع ابهام نموده و آنها را حاشیه نویسی نماید.

در این سیستم دو دسته الگوریتم جهت یادگیری ماشین بررسی شده است. در روش اول بوسیله روش استخراج قوانین پیوند که یکی از روش های داده کاوی است الگوریتمی ابتکاری جهت یادگیری سیستم طراحی شده است. این الگوریتم^{۱۰} به گونه ای تغییر داده شده تا بتواند توسط داده های متنی به استخراج قوانین پیوند میان کلمات جهت آموزش سیستم پردازد. در روش دوم مسئله انتخاب کلاس مناسب معنایی برای یک کلمه را به مسئله رده بندی^{۱۱} تبدیل نموده و بوسیله ابزار وکا^{۱۲} توسط ۱۰ الگوریتم از روش های آموزش با نظارت^{۱۳} یادگیری ماشین تست و ارزیابی شده است. در هر دو روش سیستم بوسیله داده های آموزشی حاشیه نویسی شده استاندارد مجموعه داده سنس ایوال^{۱۴} مورد تست و ارزیابی قرار گرفته و نتایج الگوریتم های مختلف مورد مقایسه قرار گرفته است. سپس راهکار های پیشنهادی برای بهبود الگوریتم های رده بندی برای رفع ابهام معنی کلمات بوسیله همین داده ها تست و ارزیابی شده و میزان بهبودی الگوریتم ها مشخص گردیده است.

۱-۴ ساختار پایان نامه

این پایان نامه دارای شش فصل به شرح زیر می باشد. پس از فصل اول، در فصل دوم به معرفی وب معنایی و فناوری های مربوط به آن مثل هستی شناسی و استنتاج از واقعی هستی و جایگاه منطق در وب معنایی پرداخته و سپس مفاهیم کلی آموزش ماشینی و متن کاوی را شرح داده ایم. در فصل سوم به کار های انجام شده قبلی در زمینه حاشیه نویسی معنایی، استخراج مفاهیم از متن پرداخته شده است. در فصل چهارم به شرح معماری و طراحی سیستم پیشنهادی پرداخته شده و فصل پنجم به تست و ارزیابی و نتایج تجربی حاصل از سیستم فوق اختصاص داده

^۹ WordNet

^{۱۰} Apriori

^{۱۱} Classification

^{۱۲} WEKA

^{۱۳} Supervised Learning

^{۱۴} SensEval-۲

شده است. در پایان در فصل ششم به نتیجه گیری و معرفی راهکار های آینده جهت بهبود این فناوری پرداخته شده است.