

## تعیین محدوده جملات فارسی به کمک درخت تصمیم گیری

سیده طاهره میرعمادیان<sup>۱</sup>، خشایار یغمایی<sup>۲</sup> و سید مرتضی سیف الله پور<sup>۳</sup>  
<sup>۱</sup> دانشجوی کارشناسی ارشد الکترونیک دانشگاه سمنان، tmiremadian@yahoo.com  
<sup>۲</sup> عضو هیات علمی دانشگاه سمنان، khyaghmaie@semnan.ac.ir  
<sup>۳</sup> دانشجوی کارشناسی ارشد الکترونیک دانشگاه سمنان، moriseif@gmail.com

### چکیده

این مقاله به بررسی تعیین محدوده جملات فارسی به کمک درخت تصمیم گیری می پردازد. تعیین محدوده جمله از جمله مراحل پیش پردازش در اکثر الگوریتمهای متن کاوی و پردازش زبان طبیعی می باشد. داده های این مقاله در حدود ۲۳۰۰۰ کلمه می باشد که به همراه برجسب هایشان از پیکره بی جن خان تهیه شده است. تعداد برجسب های این مجموعه به دلیل پرهیز از پیچیدگی زیاد، به ۱۸ برجسب تقلیل داده شده است و برجسب جدید محدوده جمله که شامل ابتدا، انتها، وسط و ابتدا/انتها می باشد به کلمات تخصیص یافته است. نتایج حاصل از شبیه سازی درخت تصمیم گیری داده های فوق به کمک نرم افزار weka، برابر ۹۶٫۹۸٪ می باشد.

### کلمات کلیدی

محدوده جملات فارسی، متن کاوی، درخت تصمیم گیری

در زبان فارسی متاسفانه کارهای زیادی صورت نگرفته است و در این زمینه تنها می توان به [7,8] اشاره نمود. در [7] تنها به رفع ابهام از علامت نشانه گذاری نقطه (".") پرداخته شده است که آیا انتهای جمله است و یا جزئی از کلمات اختصاری و مخفف می باشد. در [8] نیز به بررسی ساختاری جمله و تعیین حدود جمله به کمک چندین روش پرداخته شده است که شامل بکارگیری روش تشخیص فعل در جمله، استفاده از مدل چندتایی<sup>۶</sup> بر تعیین حدود جمله و استخراج خصیصه از جملات و به کارگیری شبکه عصبی به عنوان طبقه بندی کننده می باشد. این مقاله با هدف بهبود دقت شناسایی محدوده جملات با استفاده از استخراج تعدادی ویژگی از هر کلمه و همسایگان آن در متن و همچنین استفاده از مدل های دوتایی و به کارگیری درخت تصمیم گیری به عنوان طبقه کننده، به این امر پرداخته است که در بخش های بعدی به طور مفصل بحث خواهد شد.

### ۲- اقسام جمله از لحاظ دستوری

به لحاظ دستوری، جمله یک یا مجموع چند واژه است که بر روی هم پیام کاملی را از گوینده به شنونده می رساند. در زبان فارسی جمله ها از چند حیث دسته بندی می شوند.

### ۱- مقدمه

تعیین محدود جمله یک مسئله اساسی برای بسیاری از کاربردهای پردازش زبان طبیعی<sup>۱</sup> مانند تجزیه<sup>۲</sup>، استخراج اطلاعات<sup>۳</sup>، ترجمه ماشینی<sup>۴</sup>، خلاصه سازی<sup>۵</sup> و ... می باشد. اگرچه قطعه بندی جمله می تواند از روی علائم نشانه گذاری مانند نقطه، ویرگول و یا حتی افعال در زبان فارسی صورت گیرد، اما همچنان ابهام های زیادی در متون وجود دارد که عمل شناسایی را سخت می کند.

از کارهای انجام شده در این زمینه می توان به [1,2,3] اشاره کرد که در زبان انگلیسی کار شده است. هم چنین در [4,5] به بررسی محدوده جمله در زبان هندی پرداخته شده است و در [6] نیز زبان ژاپنی را مورد بررسی قرار داده اند. در این کارها مسئله تعیین حدود جمله تبدیل به مسئله رفع ابهام از علائم نشانه گذاری شده است که با آن به صورت یک مسئله رده بندی برخورد شده است.

<sup>۱</sup> Natural Language Processing (NLP)

<sup>۲</sup> Parsing

<sup>۳</sup> Information Extraction

<sup>۴</sup> Machine Translation

<sup>۵</sup> Abstraction

<sup>۶</sup> N-garm