

استخراج خودکار اطلاعات از متون فارسی در دامنه خاص

آرمان شریف زاده^۱، مهرنوش شمس فرد^۲

^۱ کارشناس ارشد نرم افزار، دانشکده مهندسی برق و کامپیوتر، دانشگاه آزاد اسلامی، قزوین، ایران
sharifzadeh.arman@gmail.com

^۲ آزمایشگاه پردازش زبان طبیعی، دانشکده مهندسی برق و کامپیوتر، دانشگاه شهید بهشتی، تهران، ایران
m-shams@sbu.ac.ir

چکیده

حجم انبوه متون قابل دسترس بخصوص در گستره جهانی اینترنت و اطلاعات موجود در این حجم انبوه، اهمیت استخراج خودکار اطلاعات از متن را بیشتر نشان می دهد. استخراج اطلاعات از متن شامل ارائه قالب ساخت یافته از اطلاعات دلخواه موجود در متن می باشد. در این مقاله به معرفی یک سامانه استخراج خودکار اطلاعات از متون فارسی در دامنه خاص می پردازیم. سامانه استخراج خودکار اطلاعات برای زبان فارسی در حوزه اخبار حوادث تروریستی بر اساس ترکیبی از روش های یادگیر مانند الگوریتم ماشین بردار پشتیبان و مدل میدان های تصادفی شرطی و روش های مبتنی بر الگوهای استخراج، معرفی و ارزیابی شده است. نتایج بدست آمده نشان می دهد که این سامانه در مقایسه با کارهای مشابه دارای دقت و بازخوانی قابل قبولی است.

کلمات کلیدی

استخراج اطلاعات، یادگیری ماشین، الگوریتم ماشین بردار پشتیبان، مدل میدان های تصادفی شرطی، زبان فارسی

۱- مقدمه

استخراج اطلاعات بر اساس این دو نگرش در ۳ دسته طبقه بندی می-شوند [۲]:

- (۱) رویکرد دستی با استفاده از مهندسی دانش
- (۲) رویکرد یادگیر با استفاده یادگیری ماشین
- (۳) رویکرد ترکیبی

در حقیقت سامانه های استخراج اطلاعات از رویکرد مهندسی دانش به سمت مدل های آماری در حال حرکت هستند [۱].

روش های سنتی استخراج اطلاعات عمدتاً مبتنی بر قانون و بر پایه استفاده از عبارات منظم^۱ بوده اند. اگرچه از روش های مبتنی بر گرامر^۲ هم استفاده شده است اما در مقیاس داده های حجیم و تعداد زیاد قوانین این روش دشوار به نظر می رسد. سامانه فاستاس^۳ و همچنین کار هابز^۴ و همکاران [۴] از نمونه های سنتی است که برای استخراج اطلاعات از متن ارائه شده است. برای زبان فارسی نیز یک سامانه استخراج اطلاعات با رویکرد مبتنی بر قانون و با استفاده از الگوهای استخراج با نام "مرصاد" [۵] معرفی شده است که بر روی اخبار در حوزه نظامی کار می کند. البته این گونه سامانه ها به دلیل ابهام های گسترده بخصوص در جمله های طولانی تا حدی کند و مستعد خطا هستند.

استخراج خودکار اطلاعات از موضوعات اصلی علم پردازش زبان طبیعی^۱ است که با موضوعات بسیاری مانند زبان شناسی^۲، یادگیری ماشینی^۳، ارزیابی اطلاعات^۴، پایگاه داده، وب جهان گستر و تجزیه و تحلیل اسناد در ارتباط می باشد. شناسایی خودکار نوع خاصی از موجودیت ها، روابط و یا وقایع از متن با عنوان استخراج اطلاعات از متن نامیده می شود [۱].

منابعی که در استخراج اطلاعات به عنوان ورودی مورد توجه هستند در سه نوع ساخت یافته^۵ (مثل جداول پایگاه داده یا هستان شناسی ها)، نیمه ساخت یافته^۶ (مثل فرهنگ های لغت و مستندات HTML) و غیرساخت یافته^۷ (مثل متون زبان طبیعی موجود در پیکره ها) معرفی شده اند.

جستجوی اطلاعات در ورودی غیرساخت یافته دشوارتر از جستجو در ورودی نیمه ساخت یافته و همچنین جستجو در ورودی نیمه ساخت یافته دشوارتر از جستجو در ورودی ساخت یافته است.

به طور کلی دو نگرش مختلف مهندسی دانش و یادگیری ماشین برای استخراج اطلاعات از متن وجود دارد که رویکردهای موجود برای