

بهبود دسته بندی متون فارسی در روش همسایگی وزن دار

فرزین یغمایی^۱، سعید تعبدی^۲

^۱f_yaghmaee@semnan.ac.ir

^۲s.taabbody@gmail.com

^۱ دانشگاه سمنان، دانشکده مهندسی برق و کامپیوتر

چکیده:

با رشد روز افزون منابع اطلاعاتی و حجم مقالات و مطالب تولید شده در زمینه‌های مختلف و به شکل‌های متنوع اعم از رسانه‌های مختلف دیجیتال نیاز به دسترسی آسان اطلاعات نیز افزایش می‌یابد. یکی از نیازهای اولیه در بالا بردن سرعت دسترسی به اطلاعات و پردازش این مطالب که غالباً دارای حجم بالایی نیز می‌باشند، دسته بندی^۱ این اطلاعات در طبقات مختلف می‌باشد. دسته بندی متون به عمل برچسب زدن یا تفکیک یک متن در قالب یکی از دسته‌های از پیش تعیین شده گفته می‌شود. در این مقاله به بررسی عملکرد الگوریتم^۳ WKNN با استفاده از معیار وزن‌دهی tf-idf می‌پردازیم. همچنین برای بالا بردن دقت در انتخاب طبقه صحیح و به منظور افزایش کارایی الگوریتم از روش میانگین‌گیری از داده‌ها به عنوان معیار ارزیابی استفاده می‌کنیم. نتایج به دست آمده از تفکیک متون فارسی با استفاده از روش‌های فوق نشان دهنده دقت ۸۹ درصد می‌باشد.

کلمات کلیدی:

دسته بندی متن، تفکیک کننده WKNN، نزدیکترین همسایگی، وزن‌دهی

Bayes و غیره اشاره نمود. موارد ذکر شده جزء الگوریتم‌های

پرکاربرد در حوزه^۱ IR می‌باشند [۱][۲][۳]

در زمینه دسته بندی متون انگلیسی روشها و مقالات متنوعی وجود دارد که هر کدام ویژگی‌ها خاص خود را دارا می‌باشند اما در زمینه دسته بندی متون فارسی کار چندانی صورت نگرفته است.

در این مقاله، به منظور دسته‌بندی خودکار متون فارسی از روش شاخص‌گذاری متون، تعیین میانگین وزنها و همچنین از تفکیک کننده موسوم به K مین همسایه نزدیک استفاده شده است.

یکی از راههای موثر در دسته‌بندی متون امکان ایجاد رابطه-ای بین کلمات موجود در متن و گروه مرتبط با آن متن است که در صورت یافتن چنین رابطه‌ای پیش‌بینی گروه مورد نظر امکان پذیر خواهد بود. برخی از تحقیقات انجام شده بر روی دسته بندی متون فارسی عبارتند از: منبع [۲] که در آن از ۳ روش برای طبقه بندی متون فارسی استفاده شده است و همچنین منبع [۱]

۱ - مقدمه

با پیشرفت در علوم کامپیوتری و گسترش کاربرد آن در تمام زمینه‌ها، حجم ذخیره و پردازش اطلاعات افزایش یافته و امکان دسترسی به این اطلاعات از محل‌های دیگر نیز با استفاده از شبکه‌های موجود بین کامپیوترها امکان‌پذیر شده است. به منظور بهبود نحوه دسترسی به این اطلاعات، وجود سیستم‌هایی برای انجام اعمالی از قبیل فیلترینگ، طبقه بندی، جستجو و ... ضروری می‌باشد.

دسته بندی متن فرایند تصمیم‌گیری برای نسبت دادن یک متن به یک گروه خاص از متون می‌باشد. از این سیستم برای دسته بندی متن‌ها یا سایت‌ها برای دسترسی آسان به اطلاعات آنها استفاده می‌شود. به منظور دسته بندی متون می‌توان از الگوریتم‌های متنوعی استفاده نمود که از آن جمله می‌توان به شبکه‌های عصبی، نزدیکترین همسایگی^۴، SVM^۵، Naive