

بازبینی

Two scalable algorithms for associative text
classification

داریوش حسن پور آده

۹۳۰۸۱۶۴

این مقاله در مورد طبقه‌بندی انجمنی متون^۱ می‌باشد، مساله‌ای که در مقاله بیان شده و سعی در حل این مساله داشته این است که آورده در برخی الگوریتم‌های از این نوع برای بالا بردن کارایی و دقت خود تعداد بسیار زیادی قوانین انجمنی زیادی را تولید می‌کنند، بنابراین زمان زیادی را این الگوریتم‌ها نیاز دارند و همچنین در برخی از موارد ممکن است به علت تولید بیش از حد نیاز قوانین کارایی پایینی داشته باشند. بنابراین این مقاله آمده است ۲ عدد الگوریتم همکار جهت افزایش کارایی و کاهش پیچیدگی محاسباتی برای این هدف معرفی کرده است که اولی قوانین استخراج شده را به صورت موثر ذخیره می‌کند و دیگری سرعت تطبیق قوانین^۲ بالا می‌برد.

به طور کلی هدف یک دسته‌بند متون، دسته‌بندی اسناد در قالب تعداد معینی دسته‌های از پیش تعیین شده می‌باشد. هر سند می‌تواند در یک، چند و یا هیچ دسته‌ای قرار بگیرد. این موضوع می‌تواند در قالب یک یادگیری خودکار قرار گیرد تا با استفاده از آن بتوان هر سند را به طور خودکار به دسته‌ای نسبت داد. در این مقاله، از روش دسته‌بندی بر مبنای قواعد انجمنی که از روی فرایند کاوش الگوهای مکرر مجموعه داده‌های آموزشی تولید شده استفاده کرده‌اند، می‌شود. این فرایند با فرآیندی که در داده‌های بزرگ پایگاه داده‌ها استفاده می‌شود یکسان می‌باشد. استفاده از قواعد انجمنی و ترکیب آن با قواعد دسته‌بندی و ایجاد مدل جدیدی با عنوان دسته‌بندی انجمنی و استفاده از آن برای دسته‌بندی متون می‌باشد.

در مورد نکات قوت روش ارائه شده می‌توان گفت که به علت اینکه روش ارائه شده به صورت ماژولار^۳ ارائه شده اند - یعنی به دو قسمت ذخیره‌سازی و تطبیق قوانین تقسیم شده است، کلیه فواید حاکم به سیستم‌های ماژولار را دارا می‌باشد و اینکه با توجه به اینکه تطبیق قوانین یکی از وقت‌گیرترین فاز اجرای الگوریتم‌های طبقه‌بندی انجمنی می‌باشد بنابراین با پرداختن به این فاز و تلاش برای افزایش سرعت این فاز از امتیازات این روش می‌توان به حساب آورد. از دیگر مزایای این روش می‌توان به قابلیت تفسیر ساده، درک آسان قواعد توسط انسان اشاره کرد. حذف قواعد ضعیف می‌تواند تا حد فوق‌العاده دقت دسته‌بندی را افزایش دهد. و همچنین ویژگی‌ها هم می‌توانند منفرد باشند و هم چندگانه، یعنی می‌توان از اطلاعات ترکیبی ویژگی‌های چندگانه استفاده کرد.

دسته‌بندی بر مبنای قواعد انجمنی دارای معایبی هم هست که از آن جمله افزایش بعد فضای برداری ویژگی‌ها می‌باشد که برای رفع این مشکل از تکنیک‌های کاهش بعد فضای ویژگی‌ها استفاده می‌شود، و همچنین افزایش تعداد قواعدی که در فاز آموزش تولید شده‌اند و باعث افزایش بیهوده زمان محاسبات و کاهش تاثیر در دسته‌بندی انجمنی می‌شوند. برای رفع این مشکل هم از تکنیک هرس کردن قواعد استفاده می‌شود. در این تکنیک فقط قواعدی که دارای کیفیت و تاثیر بالایی هستند انتخاب می‌شوند.

^۱ Associative Text Classification

^۲ Rule Matching

^۳ Modular