

خوشه بندی جملات فارسی مبتنی بر الگوریتم های هوش جمعی

مهدی بازقندی^۱، قمرناز تدین تبریزی^۲ و مجید وفایی جهان^۳

Mehdi_Bazghandi@ymail.com^۱

Tadayon@mshdiau.ac.ir^۲

VafaeiJahan@mshdiau.ac.ir^۳

چکیده

خوشه بندی یکی از مسائل مهمی است که امروزه بسیاری از محققین در زمینه های مختلف به آن پرداخته اند. تا کنون الگوریتم های کلاسیک زیادی در این زمینه ارائه شده است. که اغلب این روش ها دارای ناپایداری بوده و همچنین پارامترهای آن ها محدود به انتخاب کاربر می باشد. از کاربردهای خوشه بندی می توان به خوشه بندی متون و اسناد در موضوعات خلاصه سازی متون و بازیابی اطلاعات یاد کرد. در خوشه بندی جملات یک متن برای مشخص شدن جملات مشابه و نمی توان از روش مشابه آن (دسته بندی متون مشابه) استفاده کرد. بردارهایی به طول m و با مقادیر صفر بسیار زیاد پدید خواهد آمد. برای حل این مشکل، روشی جدید مبتنی بر PSO برای خوشه بندی جملات یک متن معرفی شده است. به طوریکه به جای استفاده از فاصله اقلیدسی و فاصله کسینوسی، از یک معیار جدید در محاسبه فاصله دو جمله استفاده شده است. معیاری که در آن، ارتباط معنایی کلمات با استفاده از ارتباطات آنها در متن در نظر گرفته می شود. همچنین تعیین تعداد خوشه های بهینه یکی دیگر از کارهای انجام شده در این مقاله است. برای ارزیابی یک مجموعه از خبرهای ورزشی فارسی انتخاب شده است. نتایج حاصل از ارزیابی روش پیشنهادی نشان می دهند که استفاده از خوشه بندی PSO معنایی، با تعیین تعداد خوشه های مطلوب، دقت بهتری را در خوشه بندی جملات در مقایسه با روش های دیگر، دارد.

کلمات کلیدی

خوشه بندی، الگوریتم PSO، بردارهای Context-Vector، شباهت معنایی

۱- مقدمه

یا در خوشه بندی سبد خرید مشتریان، فاصله بر اساس شباهت خرید تعیین می شود. لذا محاسبه فاصله بین دو داده در خوشه بندی بسیار مهم می باشد؛ زیرا کیفیت نتایج نهایی را دستخوش تغییر قرار خواهد داد. فاصله که همان معرف عدم تجانس است حرکت در فضای داده ها را می سازد و سبب ایجاد خوشه ها می گردد. با محاسبه فاصله بین دو داده می توان فهمید که چقدر این دو داده به هم نزدیک هستند و براین اساس در یک خوشه قرار داده می شود. توابع ریاضی مختلفی، برای محاسبه فاصله وجود دارند؛ فاصله اقلیدسی، فاصله همینگ و... با افزایش حجم منابع متنی موجود در وب چالشی که وجود داشته است آن است که چگونه کاربران می توانند به اطلاعات مورد نیاز خود در اسرع وقت و با دقت بالا دسترسی داشته باشند. خوشه بندی می تواند راه حلی برای حل این مسئله باشد. خوشه بندی می تواند نقش مهمی در انتخاب عنصرهای شایسته داشته باشد. با خوشه بندی اطلاعات می توان عنصرهای مشابه را در یک خوشه قرار داد و عنصر محوری را به عنوان نماینده خوشه انتخاب کرد. خوشه بندی در بسیاری از کاربردهای پردازش متن مانند خلاصه سازی، درک متن، ترجمه ماشینی و... استفاده می شود

پردازش داده، یکی از شاخص های بسیار مهم در دنیای اطلاعات است. خوشه بندی یکی از بهترین روش هایی است که برای کار با داده ها ارائه شده است. خوشه بندی قابلیت ورود به فضای داده و تشخیص ساختارش را امکان پذیر می نماید. لذا به عنوان یکی از ایده آل ترین مکانیزم ها، برای کار با دنیای عظیم داده ها محسوب می شود.

خوشه بندی، یافتن ساختاری در مجموعه ای از داده ها است که طبقه بندی نشده اند. به بیان دیگر می توان گفت که خوشه بندی قراردادن داده ها در گروه هایی است که اعضای هر گروه از زاویه خاصی شبیه یکدیگرند. در نتیجه شباهت بین داده های درون هر خوشه حداکثر و شباهت بین داده های درون خوشه های متفاوت حداقل می باشد. معیار شباهت در اینجا، فاصله بوده یعنی نمونه هایی که به یکدیگر نزدیک ترند در یک خوشه قرار می گیرند. به عنوان نمونه در خوشه بندی اسناد دوری و یا نزدیکی داده ها متناسب با تعداد کلمه های مشترکی که در دو سند وجود دارد و