

## خوشه بندی متون مبتنی بر مرکز دسته با استفاده از روش SVD و بهره گیری از نقاط همسایگی

محمد رضا علاقه بند<sup>۱</sup>؛ محمد رضا سعیدی محمدی<sup>۲</sup>؛ میر حسین دزفولیان<sup>۳</sup>

[mr.alagheband@basu.ac.ir](mailto:mr.alagheband@basu.ac.ir)<sup>۱</sup>

[m.saeedi@basu.ac.ir](mailto:m.saeedi@basu.ac.ir)<sup>۲</sup>

[dezfoulian@basu.ac.ir](mailto:dezfoulian@basu.ac.ir)<sup>۳</sup>

### چکیده

دسته بندی متن، یکی از شاخه هایی است که امروزه توجه گسترده ای را به خود معطوف داشته است که در آن متن های با محتویات مشابه به یک دسته اختصاص داده می شوند. در این مقاله روشی ترکیبی از خوشه بندی مبتنی بر مرکز دسته ها، K-means، به همراه روش SVD مورد مطالعه قرار گرفته که عمل خوشه بندی در آن، با معرفی دو مفهوم جدید نقاط همسایگی و ارتباط، انجام پذیرفته است. استفاده از SVD می تواند با حفظ ساختار زیربنایی در فضای ویژگی، اثر نویز را کاهش دهد. با استفاده از مفاهیم ارائه شده، روشی جهت انتخاب مراکز اولیه دسته ها و رابطه ای جهت محاسبه شباهت بردار متن ها با این مفاهیم معرفی می گردد. در انتها سه آزمایش مختلف بر روی یک مجموعه از پنجاه و شش متن انجام گرفته است. در آزمایش ابتدایی با استفاده از الگوریتم مقدماتی K-means به نتیجه ۵۷،۱۴ درصد سپس در آزمایش دوم با اعمال مفاهیم همسایگی و پیوند به نتیجه ۸۷،۵ درصد و آزمایش سوم با استفاده ترکیبی از مفاهیم ارائه شده در آزمایش دوم و الگوریتم SVD، دسته بندی متون را به ۹۲،۸۶ درصد بهبود می بخشد.

### کلمات کلیدی

خوشه بندی، K-means، SVD، نقاط همسایگی، مفهوم ارتباط

، و یا کاربرد های پزشکی در دسته بندی بافت های مغز [11] اشاره کرد. در ادامه یک روش جهت انتخاب نقاط اولیه مراکز ارائه شده تا نماینده هایی مناسب به عنوان مرکز هر خوشه انتخاب شوند. از این جهت، مفهومی به نام امتیاز نقاط کاندید را معرفی کرده ایم. برای یافتن نقاط کاندید لازم است از سه معیار استفاده شود: میزان شباهت بین دو متن، مقدار مفهوم ارتباط آنها و تعداد نقاط همسایگی آنها. در نهایت از مفهوم تجزیه مقادیر تکین (SVD) استفاده شده تا با بهره گیری از این مفاهیم خوشه بندی اسناد را بهبود بخشیم.

### ۲- خوشه بندی متون

در این بخش روند کار الگوریتم ارائه شده در این مقاله را معرفی می کنیم.

#### ۲-۱- پیش پردازش

معمولا متن ها حاوی کلمه هایی هستند که معنی خاصی ندارند. از این دست لغات می توان به حروف اضافه اشاره نمود که بهتر است از متن ها حذف گردند. بعلاوه به جای استفاده از لغاتی که در متن ها وجود دارند بایستی هنگام بررسی شباهت متن ها

### ۱- مقدمه

رشد حجم مستندات دیجیتالی در سال های اخیر، مدیریت اطلاعات را با مشکلات عدیده ای روبرو ساخته است که لازم است روش های دقیقی جهت گروه بندی این حجم از اطلاعات ارائه گردد. به همین علت تا به حال تحقیقات بسیار زیادی در زمینه خوشه بندی در دنیای بازیابی اطلاعات انجام شده است [۱،۲]. در مقاله حاضر، روش k-means [3-6] به کار گرفته شده است، در ابتدا تعداد k نقطه را به عنوان نماینده هر خوشه در نظر گرفته و سپس هر متن بسته به میزان شباهت، به یکی از این خوشه ها منتصب می شود. سپس مجددا مرکز هر خوشه محاسبه می گردد. این روند تا جایی که مراکز دسته ها تغییر نامحسوس داشته و یا اصلا تغییر نکنند تکرار میگردد.

از زمان ارائه الگوریتم K-mean تا کنون نسخه های متفاوت آن در بسیاری از سیستم ها در کاربرد های مختلف استفاده شده است. از کاربردهای آن می توان به کاربرد های تجاری چون دسته بندی هوشمند سهام در بازار بورس [7]، سیستم های پیشنهاد دهنده خرید کالا آنلاین [8]، کاربرد های صنعتی در دسته بندی انواع زغال سنگ [9]، تشخیص وضعیت ترافیک در شهر [10]