

خلاصه سازی گزینشی متون فارسی مبتنی بر خوشه بندی PSO

مهدی بازقندی^۱، قمرناز تدین تبریزی^۲، مجید وفایی جهان^۳ و علی بازقندی^۴

Mehdi_Bazghandi@ymail.com^۱

Tadayon@mshdiau.ac.ir^۲

VafaeiJahan@mshdiau.ac.ir^۳

Bazghandi@shahroodut.ac.ir^۴

چکیده

با افزایش روز افزون منابع متنی، هر روز بر گستره اطلاعات قابل دسترس برای کاربران افزوده می شود. به طوریکه دسترسی دقیق و درست به اطلاعات همواره یکی از مسائل اصلی در برخورد با آنها بوده است. سیستم های خلاصه سازی خودکار می توانند نقش مهمی در پوشش قسمت های اصلی متن و برطرف کردن محدودیت های زمان داشته باشند. در این مقاله یک سیستم خلاصه ساز مبتنی بر خوشه بندی جملات ارائه شده است. در این مقاله برای حل مسئله خوشه بندی چند روش معرفی شده است به طوریکه برای رسیدن به یک خوشه بندی با کیفیت از الگوریتم های هوش جمعی برای بهینه سازی روش های معرفی شده استفاده می شود. تمام روش های استفاده شده، ارتباط معنایی کلمات را با استفاده از ارتباطات آنها در متن، در نظر می گیرد. در نهایت بعد از خوشه بندی جملات با استفاده از معیارهای ذکر شده جملات مناسب و کاندید، از هر خوشه انتخاب می شوند.

برای ارزیابی یک مجموعه از خبرهای ورزشی فارسی انتخاب شده است. نتایج حاصل از ارزیابی روش پیشنهادی نشان می دهند که سیستم خلاصه ساز معرفی شده، خروجی های بهتری در مقایسه با روش های مشابه دیگر دارد.

کلمات کلیدی

خلاصه سازی، خوشه بندی، شباهت معنایی، کارایی، الگوریتم PSO

۱- مقدمه

سازگی است که در آن با توجه به معیارهای آماری، شهودی و یا ترکیبی از این دو تهیه می شود. از آنجا که در تولید این دسته از خلاصه ها، جملات متن تغییرات نحوی و معنایی ندارند، می توان آن را نوعی گزینش جملات قلمداد کرد. و دیگری چکیده است که تفسیری از متن اولیه است. در تولید چکیده، خلاصه مفاهیم جملات متن اصلی به شکل کوتاه تر بازنویسی می شود. فرایند خلاصه سازی به دو مدل تک سندی و چند سندی تقسیم می شود. در مدل خلاصه ساز تک سندی، ورودی سیستم خلاصه ساز تنها یک سند می باشد [1]. در خلاصه ساز چند سندی، ورودی خلاصه ساز چندین سند می باشد. پیچیدگی خلاصه ساز تک سندی به علت پیوستگی موضوع در مقایسه با چند سندی، کمتر می باشد. دسته بندی های مختلفی برای روش های خلاصه سازی ارائه شده است. در مقالات زیادی تشابه دو جمله بر اساس کلمات کاندید ان (اسم ها و فعل ها) در WordNet در نظر گرفته شده است. در [2] از یک معیار مبتنی بر اشتراک کلمات، برای تعیین شباهت جملات استفاده می کند. مشکل اصلی این روش ها در نظر نگرفتن معنای جملات می باشد. فقط به ظاهر کلمات توجه می کنند. در روشی دیگر از قانون احتمال شرطی بیزین برای احتمال حضور یک جمله در خلاصه استفاده می شود. [3]. بسیاری از مقالات از تکنیک خوشه بندی برای دسته بندی جملات مشابه استفاده می کنند. در [4] با استفاده از یک روش مبتنی بر مرکز جملات مشابه در یک خوشه قرار می گیرد و سپس با

با افزایش روز افزون منابع متنی الکترونیکی در شبکه گسترده جهانی وب، نیاز به دسترسی صحیح، سریع و آسان به اطلاعات بیشتری شود. و نیروی انسانی قابل توجهی که می تواند در کارهای دیگر به کار گرفته شود، معطوف به تسهیل بخشیدن به دسترسی به اطلاعات است. بر این اساس و بر این اساس و بر این اساس، بر این اساس و بر این اساس، خواندن کامل متون بزرگ نامناسب است. سیستم خلاصه سازی اتوماتیک این خلاء را پرمی کند. فرآیند فشرده سازی یک منبع به صورتی که حاصل، درحجم کمتر حاوی اطلاعات مهم منبع اصلی باشد را خلاصه سازی گویند. مزیت های عمده ی تولید خودکار خلاصه به وسیله ماشین عبارتند از: اندازه ی خلاصه قابل کنترل است، یعنی ماشین میتواند خلاصه را با توجه به میزان فشرده گی مورد نظر کاربر تهیه کند. محتوای آن قابل پیش بینی است. می توان مشخص کرد که هر بخش از خلاصه مربوط به کدام بخش از متن اصلی است. به طور کلی دو مدل اصلی برای خلاصه سازی متن وجود دارد. خلاصه گزینشی که پرکاربردترین نوع خلاصه