



ارائه یک ابزار ارزیابی خودکار خلاصه‌سازهای چکیده‌ای فارسی با بهره‌گیری از شبکه واژگان

احمد استیری، محسن کاهانی، آصف پورمعصومی و محسن عباسی

ahmad.estiri@stu.um.ac.ir, kahani@um.ac.ir, {as.poormasoomi, abasi.mohsen}@stu.um.ac.ir

آزمایشگاه فناوری وب، دانشگاه فردوسی مشهد

چکیده

امروزه با رشد چشمگیر اسناد منتشر شده در وب و نیاز اساسی به نگهداری، دسته‌بندی، بازیابی و پردازش آنها، توجه به پردازش زبان طبیعی و بهره‌گیری از ابزارهایی نظری خلاصه‌سازهای خودکار و مترجم‌های ماشینی، بیش از پیش احساس می‌شود. در این مقاله، ابزاری به منظور ارزیابی سیستم‌های خلاصه‌سازی چکیده‌ای خودکار ارائه شده است که البته از آن در دیگر کاربردهای پردازش زبان طبیعی و بازیابی اطلاعات هم می‌توان استفاده نمود.

این ابزار شامل معیارهایی برای تعیین کیفیت خلاصه‌ها به صورت خودکار از طریق مقایسه آنها با خلاصه‌های تولید شده توسط انسان (خلاصه‌های ایده‌آل) می‌باشد. بدیهی است برای انجام مقایسه متن در سطح معنا در مورد خلاصه‌های چکیده‌ای، مقایسه‌ی ظاهر لغات کافی نمی‌باشد و بهره‌گیری از شبکه‌ی واژگان، ضروری به نظر می‌رسد که با ایده‌ای مناسب برای زبان فارسی به کار گرفته شده و نتایج حاصل از ارزیابی را به طور قابل توجهی بهبود بخشیده است. این ابزار پس از طراحی و پیاده‌سازی، توسط افراد خبره در زمینه زبان‌شناسی بررسی گردید. نتایج حاصل از این بررسی‌ها قابل توجه می‌باشد.

کلمات کلیدی

پردازش زبان طبیعی، زبان فارسی، معنگرایی، ارزیابی، خلاصه‌سازهای ماشینی، شبکه واژگان

نتایج حاصل از مقایسه خلاصه ماشینی با هر یک از خلاصه‌های انسانی و یا بیشینه‌ی امتیاز حاصل شده باشد.

قبل از مقایسه متن خلاصه، جهت استانداردسازی متن بازیستی پیش‌پردازش‌هایی روی آنها انجام شود. طبیعتاً هر چه این پیش‌پردازش‌ها قوی‌تر باشد، نتایج حاصل از مقایسه متن، قابل اطمینان‌تر خواهد بود. لازم به ذکر است از آنجایی که زبان فارسی جزو زبان‌های غیر ساختیافته می‌باشد، پیچیدگی‌ها و مشکلات بیشتری دارد.

در روند مقایسه متن خلاصه، از یک ابزار ریشه‌یاب جهت محاسبه‌ی ریشه‌ی تک‌تک لغات موجود در متن استفاده می‌شود و در روند ارزیابی، ریشه‌ی لغات با یکدیگر مقایسه می‌گردد. اما بدیهی است که صرفاً مقایسه‌ی ظاهر لغات با یکدیگر برای ارزیابی خلاصه‌های چکیده‌ای، نتیجه‌ی مطلوبی نخواهد داشت؛ چرا که در خلاصه‌های چکیده‌ای، لغات تغییر می‌یابند و جملات خلاصه می‌توانند بازتولید گردد، بنابراین مقایسه با حالتی که خلاصه‌ها گرینشی یا انتخابی باشند، متفاوت خواهد بود. به همین جهت، بهره‌گیری از یک شبکه‌ی واژگان، ضروری به نظر می‌رسد.

۱ - مقدمه

در وب امروزی با توجه به گسترش روزافزون حجم اطلاعات و داده‌های انتشاریافته، دسترسی درست و مطالعه اطلاعات مورد نیاز در کوتاه‌ترین زمان ممکن، همواره یکی از چالش‌های محققان و پژوهشگران حال حاضر می‌باشد.

خلاصه‌سازی خودکار سند، یعنی تولید یک نسخه مختصرتر از سند اصلی توسط یک برنامه کامپیوتری به نحوی که ویژگی‌ها و نکات اصلی سند اولیه حفظ شود [1].

با توجه به اهمیت بسیار زیاد خلاصه‌سازها و گستردگی روش‌های ارائه شده برای خلاصه‌سازی خودکار متن، ارزیابی دقیق و صحیح این روش‌ها از اهمیت خاصی برخوردار می‌باشد. دو رهیافت در ارزیابی سیستم‌های خلاصه‌ساز وجود دارد: قضاویت انسانی و مقایسه با خلاصه‌ی مرجع. به منظور ارزیابی خلاصه‌ای که یک ماشین از یک متن تولید می‌کند، می‌توان آن خلاصه را با خلاصه‌های تولید شده توسط انسان‌ها مقایسه کرد. برای کاهش تاثیر سلیقه‌ها و نظرات شخصی در خلاصه‌های انسانی، هر خلاصه‌ی ماشینی با چند خلاصه انسانی متفاوت از همان متن مقایسه می‌گردد و نتیجه‌ی نهایی می‌تواند میانگین