

## مروری بر روش های خلاصه سازی خودکار متون

لیلا طالب علی<sup>۱</sup>، نوشین ریاحی<sup>۲</sup>

۱- گروه مهندسی کامپیوتر، دانشگاه الزهرا (س)، تهران

۲- گروه مهندسی کامپیوتر، دانشگاه الزهرا (س)، تهران

### چکیده

امروزه پردازش زبان طبیعی در زمینه های گوناگون نظیر خلاصه سازی خودکار و مترجم های ماشینی<sup>۱</sup>، توجه زیادی را به خود جلب نموده اند. در زبان فارسی هم مانند سایر زبانهای دیگر دنیا تلاش هایی در زمینه ساخت ابزارهای خلاصه سازی صورت گرفته است. تمرکز محققان بر ارایه روشهایی متمرکز است که بتواند خلاصه هایی پر محتوا، سلیس و روان نسبت به روشهای خلاصه سازی پیشین ارایه دهند. خلاصه سازی یک مهارت نگارشی به شمار می رود، که هدف از ایجاد سیستم خلاصه ساز اتوماتیک تقلید کلیه مراحل است که توسط عامل انسانی انجام می شود، بدین صورت که متن به طور کامل خوانده و فهمیده شود و با تشخیص و تفکیک قسمت های مهم و غیرمهم متن، نسخه خلاصه شده متن اصلی تولید گردد.

هدف از خلاصه سازی خودکار سند، تولید یک نسخه مختصرتر از سند اصلی توسط یک برنامه رایانه ای به نحوی که ویژگی ها و نکات اصلی سند اولیه حفظ شود. بنابر تعریف ارائه شده در استاندارد ISO ۲۱۵، خلاصه "یک بازگویی مختصر از سند" می باشد. روش های تولید خلاصه را با توجه به انواع دیدگاه های مختلف خلاصه سازی می توان به چندین دسته تقسیم بندی نمود، به عنوان مثال روش تولید خلاصه تک سندی و چند سندی، تک زبانه و چند زبانه، مبتنی بر تعامل با کاربر و غیر مبتنی بر تعامل با کاربر و... لیکن با توجه به اهمیت فاکتور خروجی در تولید خلاصه در این مقاله به بررسی روش های تولید خلاصه گزینشی (استخراجی) و چکیده ای (مفهومی) پرداخته می شود.

**واژگان کلیدی:** پردازش زبان طبیعی، خلاصه سازیهای ماشینی، روابط کلامی، تشابه معنایی، شبکه

واژگان

<sup>۱</sup> . Machine Translator

## مقدمه

خلاصه سازی خودکار متن به عنوان هسته مرکزی طیف گسترده ای از ابزارهای پردازش گر متن مانند خلاصه سازی ماشین، سیستم های تصمیم یار، سیستم های پاسخ گو، موتورهای جستجو، از سال ها پیش مطرح شده و همواره به عنوان یک موضوع مهم مورد بررسی و تحقیق قرار گرفته است. (Mekeown et al, ۲۰۰۳) آغاز فعالیت سیستم های خلاصه سازی خودکار متن مربوط به سال ۱۹۵۰ می شود. به دلیل کمبود کامپیوترهای قدرتمند و مشکلات موجود برای پردازش زبانهای طبیعی (NLP)، کارهای اولیه بر روی مطالعه ظواهر متن مانند موقعیت جمله و عبارات اشاره متمرکز شده بود. سال ۱۹۷۰ تا ۱۹۸۰ هوش مصنوعی بکار آمد (Buyukkokten et al, ۲۰۰۳)، (شهاب، ۱۳۸۱) و (Mazdak, N., ۲۰۰۴). Hassel, ایده ی AI، استخراج نمایش های دانش مانند فریمها یا الگوها برای شناسایی موجودیتهای مفهومی از متن و استخراج روابط بین موجودیت ها با مکانیزمهای استنتاج بود. از اوایل ۱۹۹۰ تا به حال هم روشهای بازیابی اطلاعات (IR) بکار گرفته شده است (Dalianis, ۲۰۰۰)، (کریمی و شمس فرد، ۱۳۸۵) و (مشکی و آنالویی، ۱۳۸۸). بیشتر این روش ها بر روی سطح سمبولیک متمرکز بوده و وارد حوزه های معنایی نمی شوند.

Kupiec اولین الگوریتم را در این زمینه پیشنهاد داد (کریمی و شمس فرد، ۱۳۸۵). در این روش بر اساس مقادیر ویژگیهای یک جمله، احتمال حضور آن در خلاصه تخمین زده می شود. او عمل خلاصه سازی را به صورت یک مسئله دسته بندی، در نظر گرفت و دسته بندی کننده های بیزین را برای تعیین جملاتی که باید در خلاصه وارد شوند بکار برد. Chuang و Yang چندین الگوریتم مانند درخت تصمیم و دسته بندی کننده را برای استخراج قطعات جمله پیشنهاد داد. (اخوان و همکاران، ۱۳۸۷) در سال ۱۹۹۷، Barzily روشی برای تولید خلاصه با پیدا کردن زنجیره های لغوی معرفی کرد که به توزیع کلمه و اتصالات لغوی بین آنها برای تقریب زدن محتوا و ارائه یک نمایش از ساختار لغوی بهم پیوسته متن اتکا می کرد. (بهره پور و همکاران، ۱۳۸۷) از روش های آماری هم در خلاصه سازی متن بسیار استفاده شده است که از جمله آنها می توان به روش های مبتنی بر مدل موضوع و روش های مبتنی بر گراف اشاره نمود.

از اوایل سال ۲۰۰۰ به تدریج بحث خلاصه سازی مبتنی بر کاربر و یا خلاصه سازی شخصی سازی شده مطرح شد و تا به امروز نیز مقالات بسیاری در این زمینه منتشر شده است. ایده اصلی خلاصه سازی شخصی سازی شده و یا مبتنی بر کاربر این است که کاربران مختلف با توجه به دانش و پس زمینه اطلاعاتی که دارند، دیدگاههای متفاوتی روی اسناد یکسان دارند.

علی رغم اینکه در زمینه خلاصه سازی متن کارهای زیادی از سالهای دور انجام شده است ولی همچنان شاهد مقالات متعددی هستیم که همه ساله در این زمینه منتشر می شوند و سعی در بهبود دقت روش های پیشین دارند. (R. Ferreira, ۲۰۱۴).

## جنبه ها و مدل های مختلف خلاصه سازی متن

روش های خلاصه سازی متن معمولاً از دیدگاه های مختلفی تقسیم بندی می شوند. این دسته بندی را در چند قالب اصلی می توان تشریح کرد (ستوده و همکاران، ۱۳۸۹).

- بر اساس منبع ورودی: بر اساس نوع منبع ورودی خلاصه سازها به دسته های کلی تک سندی و چند سندی تقسیم بندی می شوند. همچنین بسته به این که نوع متن، اخبار، علمی، گزارش، عمومی و... باشد، روش برخورد هم متفاوت می باشد. به عنوان مثال، در خلاصه سازهای اخبار یکی از ساده ترین روش هایی که استفاده می شود این است که تیتز متون به همراه جمله های اول پاراگراف به عنوان خلاصه استفاده می شود و این در حالی است که این روش برای متون علمی جوابگو نمی باشد. همچنین بر اساس زبان متن هم دو دسته خلاصه ساز تک زبانه و چند زبانه داریم. روش خلاصه سازی بسته به اینکه متون تک زبانه و یا چند زبانه باشد، متفاوت می باشد. چرا که زبان های مختلف با یکدیگر تفاوت داشته و بسیاری از ویژگی ها که در یک زبان صادق است در زبان دیگر ممکن است صادق نباشد.
- بر اساس هدف: در این تقسیم بندی باید مشخص شود که متن خلاصه به چه منظوری تولید می شود. هدف خلاصه سازی می تواند هشدار، پیش نمایش، آگاهی، زندگی نامه و... باشد و این موضوع در تعیین روش خلاصه سازی تاثیر گذار می باشد. همچنین خلاصه ها می توانند عمومی یا مبتنی بر پرس و جو باشند و یا اینکه مبتنی بر کاربر و یا کلی باشند.
- بر اساس خروجی: بر این اساس خلاصه ها به دو دسته عمده خلاصه سازهای چکیده ای و استخراجی تقسیم بندی می شوند.

### دسته بندی الگوریتم های خلاصه سازی

به طور کلی چهار دسته روش برای روند خلاصه سازی وجود دارد که عبارتند از:

۱. روش های آماری سطحی: خلاصه سازهای سطحی روش های آماری، مکانی و تحلیل های مکانی - آماری انجام می دهند. این روشها از یک سری خصوصیات ساده متن برای نمایش اطلاعات محتوا استفاده می نمایند. از جمله این خصوصیات می توان از خصوصیات کمی متن، فرکانس عبارت، خصوصیات محلی مانند: مکان جمله در متن، مکان محلی پاراگراف، خصوصیات زمینه ای مانند: وجود کلمات عنوان، کلمات و عبارات خاص نام برد. (عبارات خاص را می توان خلاصه هایی که در خود متن ورودی موجود است دانست).
۲. روش های درک متن: تفاوت ایده اصلی این روش با ایده روش سطحی آن است که یک نمایش داخلی از متن ورودی با مدل کردن متن و ارتباطات بین آن ها ایجاد می کند. ارتباط بین موجودیت ها شامل شباهت، هم معنی، ارتباط نحوی و ارتباط معنایی مانند: تضاد، سازگاری متن، ارتباط لغوی و هم مکانی کلمات می باشد. روش های بر اساس زبان شناسی با تجزیه کلمه و مشخص کردن نحو کلمه آغاز می شود و از پایگاه داده لغوی استفاده می نماید تا ارتباط بین جملات را استخراج نماید و جملات مهم بر اساس معنایشان انتخاب می شوند.
۳. روش های مبتنی بر یادگیری نظارتی: یادگیری نظارتی گونه ای از یادگیری است که در آن به عامل یادگیرنده تعدادی جنبه و مقدار متناظر برای خصوصیتی که بایستی یاد گرفته شود، داده شده و از او خواسته می شود تا در موارد مشابه بر اساس آنچه در ابتدا به او داده شده است عمل نماید. در این روش، با مجموعه های از سندهای آموزش و خلاصه های استخراجی آن ها، فرآیند خلاصه سازی به عنوان یک مسئله طبقه بندی مدل می شود: جملات به "جملات خلاصه" و "جملات غیر خلاصه" بر اساس مشخصه هایی که دارا هستند، طبقه بندی می شوند.

۴. روش های ساختار کلامی: روش هایی از ساختار کلامی برای تشخیص بخش های مهم متن استفاده می نمایند، این سطح پردازش متن ساختار عمومی متن و ارتباط آن ها را مدل می نماید. این ساختار شامل هدف سند، موضوعاتی که در متن آمده است و همین طور شامل ساختار کلامی متن می باشد. این مسئله که چگونه بتوان یک ساختار را به صورت اتوماتیک استخراج نمود و چگونگی استفاده از آن در تبدیل نحوی اطلاعات مسئله مهمی است.

## روش های آماری سطحی

مزیت اصلی این گونه روش ها سادگی آن ها و همچنین امکان اعمال بر روی هر نوع متن ورودی می باشد. همچنین از آنجایی که هیچ تلاشی در آن برای فهم متن صورت نمی پذیرد، کارایی و رواج بالایی دارند. از آنجایی که در این دسته از روش ها ارتباط بین کلمات در نظر گرفته نمی شود، عدم انسجام و پیوستگی در خلاصه های تولیدی به وضوح ملاحظه می شود، که می توان این مشکل را عیب مهم این دسته از روش ها دانست.

- ویژگی های ادمونسونی<sup>۱</sup>: احتمالاً تاریخیچه اولین گروه از روش های خلاصه سازی به زمان ادمونسون باز می گردد. زیرا بسیاری از روش هایی که پس از آن تولید شدند، تلاش نمودند جنبه های گوناگون این سیستم و نحوه تأثیر آن ها را به گونه ای مختلف تغییر دهند. برخی از این ویژگی ها جهت امتیازدهی به جملات مورد استفاده قرار می گیرد، که عبارتند از (بازقندی و تدین تبریزی، ۱۳۹۱):

- فراوانی واژه ها: مهم ترین و پرکاربردترین جنبه در میان ویژگی های ادمونسونی برای امتیازدهی، میزان تکرار کلمات گوناگون موجود در متن یا متون ورودی است. بکارگیری این جنبه بر پایه فرض واژه های موضوعی<sup>۲</sup> است. این فرض بیان می دارد که واژه هایی که به نسبت در یک متن بیشتر بکار رفته اند، دارای بار معنایی بیشتری هستند. از آنجایی که در یک متن با موضوعی خاص، واژه های زیادی مرتبط با آن موضوع بکار می روند، این فرض منطقی به نظر می رسد. البته بکارگیری این ایده به صورت خام کارایی چندانی نخواهد داشت؛ زیرا واژه هایی مانند حروف اضافه، حروف ربطی و ضمائر اشاره در هر متنی به صورت پرتکرار وجود دارند، اما بار معنایی مهمی ندارند. یک ایده برای جلوگیری از اثر مخرب کلمات پرکاربرد و کم اهمیت این است که میزان اهمیت یک کلمه را در یک متن، علاوه بر تعداد تکرارهای آن، به صورت معکوس به تعداد اسنادی که در یک پیکره دارای این واژه می باشند، وابسته نمود. به این نحوه امتیاز دهی به هر کلمه،  $TF/IDF$ <sup>۳</sup> گفته می شود، که معادل فرکانس کلمه - معکوس فرکانس سند می باشد، و برای هر واژه به صورت زیر محاسبه می گردد:

$$tf.idf(term) = tf(term). \log\left(\frac{NumDoc}{NumDoc(term)}\right) \quad (1)$$

که در آن  $NumDoc$  تعداد کل اسناد موجود در پیکره،  $NumDoc(term)$ ، تعداد اسنادی که دارای واژه  $term$  هستند و  $tf(term)$ ، تعداد تکرارهای واژه  $term$  در سند یا اسناد ورودی می باشد. حال در صورتی که

<sup>۱</sup>. Edmonsonian Features

<sup>۲</sup>. Term Frequency

<sup>۳</sup>. Thematic Term Assumption

<sup>۴</sup>. Term Frequency- Inverse Document Frequency

متن ورودی تنها حاوی یک سند به جای چندین سند باشد، این معیار به صورت فرکانس کلمه - معکوس فرکانس جمله یا  $TF\_ISF(w,s)$  بیان می گردد. فرکانس جمله تعداد جملات سند است که حاوی کلمات هستند. این مشخصه برای تمام کلمات هر جمله محاسبه می شود (Barzilay et al, ۲۰۰۲).

- مکان واحدهای متنی: در (Chuang and Yang, ۲۰۰۰) نشان داده شده است که جملات مهم در ۸۵ درصد موارد در ابتدای پاراگراف و در ۷ درصد موارد در انتهای آن قرار دارند. همچنین یک سیاست بهینه برای مکان یابی جملات مهم ارائه شده است. بر اساس مطالعه ای در ۱۳۰۰۰ خبر از مجموعه Ziff-Davis در رابطه با محصولات کامپیوتری، عنوان هر مقاله و پس از آن جمله اول از هر پاراگراف دارای بیشترین بار معنایی می باشند. از سوی دیگر در خبرهای مجله وال استریت، پس از عنوان جملات موجود در پاراگراف اول دارای بار معنایی بیشتری می باشند. نتایج این تحقیق نشان می دهد که مکان جملات مهم بسته به منبع آن جمله، متفاوت است، که این خود یک چالش برای این مبحث می باشد.

- عبارات توضیحی: این عبارات ها وابسته به یک دامنه خاص می باشند و بر اساس تناوب تکرار از یک پیکره استخراج می شوند. این عبارات ها را می توان به دو دسته عبارات تنبیهی و تشویقی تقسیم نمود. عبارت تشویقی احتمال آمدن واحدهای متنی پر ارزش را در خلاصه افزایش می دهد و عبارات های تنبیهی این احتمال را کاهش می دهد. به عنوان مثال عبارت "نتیجه می گیریم که" یک عبارت تشویقی و عبارت "وی گفت" یک عبارت تنبیهی می باشد.

- امضاء سرفصل: این روش نیز استراتژی دیگر برای یافتن جملات پر اهمیت موجود در متن و امتیازدهی به آنهاست. در این روش آماری بدون استفاده از پایگاه های دانش و یا پارسرها جملات پر اهمیت استخراج می گردد. ایده اصلی این روش این گونه است که واژگانی که از نظر معنایی با یکدیگر مشابهت دارد در محیط های یکسانی ظاهر می شوند. به عبارت دیگر این روش بیان می دارد که اگر یک ایده یا اتفاق خاص مهم و کلیدی در یک متن وجود داشته باشد، مجموعه هایی از واژگان با ساختاری قابل پیش بینی در متن وجود خواهد داشت که یک موضوع را به نحوی یکتا شرح می دهند. امضاء سرفصل با معرفی یک زوج متشکل از یک مبحث (موضوع) و مجموعه ای از واژگان مرتبط با آن، این ساختارها را به صورتی رسمی بیان می کند. برای یافتن این امضاءهای سرفصل معمولاً از روش های باهم آیی استفاده می شود. در (هنر پیشه، ۱۳۸۶) مدلی بهبود یافته از امضاءهای سرفصل ارائه شده است. در این مدل علاوه بر پیدا نمودن امضاءهای سرفصل، ارتباط میان واژه های موجود در امضاءهای سرفصل نیز استخراج می شود. این ارتباطات یا به صورت ارتباطات نحوی و یا به صورت باهم آیی می باشد.

- روش باهم آیی: باهم آیی کلمات بر اساس این ایده عمل می کند که کلماتی که در یک زمینه مشترک با هم مشاهده می -شوند، با یکدیگر در ارتباط باشند. چند- وزنی ترتیبی از چند کلمه است. چند- وزنی زمانی در یک سند وجود دارد که این کلمات به یک ترتیب پشت سرهم در سند ظاهر شود. در اکثر مراجع این توالی کلمات را با در نظر گرفتن ریشه کلمات بدون در نظر گرفتن کلمه های توقف محاسبه می نمایند.

## روش های درک متن

این گونه روش ها، روش های پرهزینه ای هستند زیرا علاوه بر نیاز به منبع دانش، نیازمند تفسیر متن پیچیده می باشد. ولی عملکرد بهتری نسبت به روش اول دارند. از آنجایی که ایده اصلی این نوع خلاصه سازی ها بر پایه درک متن می باشد نیاز به خط مشی های بر پایه دانش و مقایسه بخش هایی از متن با پایگاه داده موجود دارد. تحلیل نحوی و معنایی دو بخش مهم این نوع خلاصه سازی می باشد.

یک نمونه از اطلاعاتی که در بسیاری از کاربردهای متن کاوی اهمیت فراوان دارد، وابستگی بین اجزای متن است. دو دسته مهم از وابستگی های متنی عبارتند از: همبستگی<sup>۱</sup> و ارتباط معنایی<sup>۲</sup>. پدیده همبستگی، معادل با این واقعیت است که بعضی از عناصر متنی مانند واژه ها تمایل دارند که در کنار هم ظاهر شوند، در حالی که پدیده ارتباط معنایی بر این حقیقت اشاره دارد که یک ارتباط هوشمندانه بین جملات متن وجود دارد. ارتباط معنایی یک ارتباط سطح بالاتر از همبستگی است. هر دو پدیده همبستگی و ارتباط معنایی پدیده هایی ذهنی<sup>۳</sup> هستند، به همین دلیل شناسایی آن ها با چالش هایی مواجه است. شناسایی همبستگی ساده تر از ارتباط معنایی است، زیرا با بررسی فراوانی واژه ها و وقوع هم زمان آن ها قابل شناسایی است (ارژنگ، ۱۳۸۳).

### • روش مبتنی بر گراف (R. Ferreira, ۲۰۱۴):

در این گونه روش ها بعد از گام های پیش پردازش برای مشخص کردن عناوین (اسامی مهم و...) در متن، سند به شکل گرافی غیر جهت دار که جملات گره های تشکیل دهنده آن هستند ارائه می شود و تئوری گراف می تواند به راحتی برای تجسم شباهت بین سندی و درون سندی به کار رود (Hovy and Marcu, ۱۹۹۸). یکی از الگوریتم هایی که برای پیدا کردن جملات با ارزش مورد استفاده قرار می گیرد الگوریتم TextRank می باشد. الگوریتم Page Rank اولین بار در سال ۱۹۹۲ توسط Larry Page و Sergey Brin در دانشگاه استنفورد ارائه شده است، این الگوریتم یک روش مستقل از پرس و جو می باشد. الگوریتم رتبه بندی صفحات مهم ترین مهمترین عامل در اندازه گیری اعتبار یک سایت می باشد. امروزه در دنیای اقتصادی وب سایت هایی که رتبه PR بالاتری دارند با مبالغ هنگفتی معامله می شوند و همچنین برای موتور های جستجو عامل مهم برای دادن پاسخ به درخواست کاربر ها همین معیار رتبه بندی است. این روش یک بار به هر سند وب امتیاز اختصاص داده و از این امتیاز، با در نظر گرفتن معیاری با توجه به پرس و جوی کاربر جهت رتبه بندی اسناد استفاده می کند. این الگوریتم رتبه هر صفحه را با اختصاص وزن به پیوندی که به آن صفحه داده شده است به دست می آورد که مقدار این وزن به کیفیت صفحه هایی که پیوند در آن قرار گرفته، بستگی دارد. در این صورت پیوندهای صفحات مهمتر وزن بیشتری میگیرند. جهت مشخص کردن کیفیت صفحه های رجوع کننده، در Page Rank از رتبه آن صفحه که به صورت بازگشتی تعیین و مقدار اولیه آن اختیاری است، استفاده می شود.

اگر  $n$  سند در دسترس باشد، مقدار اولیه رتبه سند را می توان برابر  $1/n$  در نظر گرفت.

$$PR(A) = (1 - d) + d \cdot \sum \frac{PR(T_i)}{C(T_i)} \quad (2)$$

<sup>۱</sup>. Cohesion

<sup>۲</sup>. Coherence

<sup>۳</sup>. Subjective



فرمول فوق جهت محاسبه رتبه یک صفحه مانند  $A$  استفاده می شود. در این فرمول،  $PR(T_i)$  رتبه صفحه  $T_i$ ،  $C(T_i)$  تعداد لینکهای خروجی صفحه  $T_i$  و  $d$  ضریب اثر میباشد که دارای مقداری بین ۰ و ۱ است (فتوحی، ۱۳۹۲). الگوریتم TextRank در واقع کاربرد الگوریتم PageRank در زمینه پردازش زبان طبیعی است. به وسیله این الگوریتم می توان کلمات مهم یا همان کلمات کلیدی را استخراج نمود. همچنین می توان جملات مهم را بر اساس میزان ارتباط با سایر جملات یا بقیه جملات با آن تعیین نمود.

## • روش زنجیره لغوی :

زنجیره های لغوی به عنوان خوشه هایی از کلمات از نظر معنایی مرتبط با هم تعریف شده اند. به عنوان مثال "خانه، سرا، اتاق" یک زنجیره است، در حالی که خانه و سرا هم معنی هستند ولی اتاق بخشی از خانه است. یک زنجیره لغوی<sup>۱</sup> مجموعه ای از لغات مرتبط با یکدیگر است که واحدی موضوعی از متن را پوشش می دهد. (ریاحی و عزالی، ۱۳۹۳) به عبارتی دیگر اگر زیر بخش های موضوعی در ارتباط با یک مطلب خاص در متنی وجود داشته باشد، یک زنجیره لغوی مرتبط، آن موضوع را متمایز می سازد. نخستین روش محاسباتی استخراج زنجیره های لغوی، برای استفاده در سیستم های خلاصه ساز ارائه گردید (حافظی و همکاران، ۱۳۸۵). این روش با بکار گیری مجموعه ای از ارتباطات میان لغات، مانند تکرار، مترادف<sup>۲</sup>، زیر مجموعه<sup>۳</sup> (به عنوان مثال زیر مجموعه بودن پراید در گروه ماشین ها)، تضاد<sup>۴</sup> و زیر بخشی (به عنوان مثال زیر بخش بودن برگ برای درخت)، زنجیره های لغوی را تولید می نماید. نحوه ساخت زنجیره لغوی به این صورت است که ابتدا زیر مجموعه ای از لغات پرتکرار موجود در متن انتخاب می شوند. به ازای هریک از لغات، اگر اتباطی با یکی از زنجیره ها موجود بود، آن لغت با ارتباطش به زنجیره متصل می شود؛ در غیر این صورت برای آن لغت یک زنجیره مجزا تشکیل داده می شود. اگر یک لغت دارای چندین معنا باشد برای هریک از معنای آن گروه زنجیره ای متفاوتی تشکیل می دهد. از میان گروه زنجیره های متفاوت گروه زنجیره ای انتخاب می شود که دارای ارتباطات بیشتری میان لغت هایش باشد. در پایان بر اساس گروه زنجیره انتخابی به هر یک از جملات امتیازی نسبت داده می شود. نحوه امتیازدهی به جملات بر اساس لغات موجود در آن هاست؛ به این ترتیب که به هر لغت به اندازه تعداد ارتباط هایش با دیگر لغات در زنجیره امتیاز داده می شود.

مشکل مهم محاسبه زنجیره ها با استفاده از شبکه لغات، تعداد زیاد معنای کلمات در زنجیره لغات است که باعث رشد نمایی تعداد گروه زنجیره ها می گردد. در صورتی که تعداد گروه زنجیره ها از یک حد آستانه بزرگ تر شود، در هر مرحله گروه زنجیره هایی انتخاب می شوند که امتیاز بالاتری داشته باشد. دو الگوریتم مختلف برای ایجاد زنجیره ها وجود دارد: رویه غیر مبهم حریصانه که زنجیره یک کلمه فقط به وسیله زنجیره های کلمات قبل

<sup>۱</sup> . Lexical Chain

<sup>۲</sup> . Synonymy

<sup>۳</sup> . Hypernymy

<sup>۴</sup> . Antonymy

از آن در متن مشخص می شود به این صورت که بر اساس روابط تعیین شده اگر زنجیره مرتبط با کلمه در زنجیره- های از قبل موجود یافت شود، آن را در همان زنجیره درج می کند و در غیر این صورت برای آن زنجیره جدید ساخته می گردد. در مقابل الگوریتم غیرحریصانه تا هنگامی که همه کلمات در متن پردازش شود منتظر می ماند سپس با توجه به تمام لغات متن زنجیره هر کدام را یافته یا ایجاد می کند. به طور کلی اکثر خلاصه سازیهای بر اساس زنجیره لغوی روش یکسانی را در مرحله پردازش متن به صورت ذیل دنبال می کنند:

الف- تولید زنجیره های لغوی

ب- امتیاز دادن به زنجیره ها

ج- یافتن قویترین این زنجیره ها جهت ایجاد خلاصه متن  
روش استفاده از رویدادها<sup>۱</sup>:

اخیراً رویدادها به عنوان یکی از راه های مدل سازی متن مورد توجه قرار گرفته اند. رویداد در هر حوزه تعریف مخصوص به خود را دارد. زبان شناسی به تعریف معنایی از رویداد و تعریف معنایی ساختار افعال می پردازد. بازیابی اطلاعات در بخش یافتن و ردیابی موضوع به رویدادها به عنوان زیر موضوعاتی نگاه می کند که می تواند از آن ها برای دسته بندی استفاده نماید. در این تعریف رویدادها می توانند رخدادها و دلایل آن ها، نتایج و هر گونه تأثیر اتفاق ها را شامل شوند. استخراج اطلاعات رویدادها را قالب هایی از پیش تعیین شده و ساخت یافته در نظر می گیرد که یک عمل را به انجام دهندگان، زمان و دیگر موجودیت های شرکت کننده مرتبط می سازد. این تعریف همان تعریف اتفاق های ساده در زبان شناسی می باشد. به طور کلی اتفاق ها مجموعه ای از فعالیت ها همراه با موجودیت های مرتبط می باشد.

در این روش ابتدا می بایست مجموعه موجودیت های مهم (مانند: فرد، زمان و مکان ها) را که رابطه نامیده می شوند، استخراج نمود، سپس ارتباط میان این رابطه ها را که همان فعل ها و اسم های فاعلانه می باشند و به آن ها ارتباط دهنده گفته می شود، استخراج کرد. برای امتیاز دهی به جملات در این روش از میزان باهم آیی رابطه ها استفاده می شود (حسامی فردو همکاران، ۱۳۸۴).

## روش های مبتنی بر یادگیری نظارتی

تاکنون روش های یادگیری نظارتی و بدون نظارتی گوناگونی برای بدست آوردن یا تقریب زدن مدلی که انسان ها برای خلاصه سازی استفاده می نمایند، ارائه شده است. اما در رابطه با این که به طور کلی کدامیک از روش ها بهتر عمل می نماید، اختلاف نظر وجود دارد. از دلایل این امر می توان به اختلاف نظر در اهمیت موضوعات گوناگون در متون و هزینه زیاد تولید یک پیکره با اندازه مناسب اشاره نمود (Porter, ۱۹۸۰).

برای تعیین میزان اهمیت واحدهای متنی، می توانیم هریک از جنبه های مطرح پیشین برای امتیازدهی را بکار ببریم. این که امتیاز ارائه شده توسط هر یک از روش ها چه مقدار در امتیاز نهایی آن ها می تواند تأثیر داشته باشد، توسط روش های

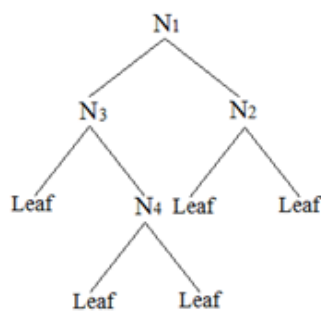
<sup>۱</sup>. Event



یادگیری ماشین تعیین می شود. دو نمونه از الگوریتم های یادگیری نظارتی را که در خلاصه سازی بکار گرفته شده اند، مورد بررسی قرار خواهد گرفت.

## • روش های مبتنی بر درخت های تصمیم<sup>۱</sup>:

در (عربی نرئی و همکاران، ۱۳۸۶) یک سیستم خلاصه سازی بر پایه الگوریتم های  $C4.5$  برای استخراج قوانین با استفاده از درخت های تصمیم ارائه شده است. البته برای بکار گیری این الگوریتم در خلاصه سازی چند تغییر در آن اعمال شده است. یکی از این تغییرات، استفاده از یک روش بهینه سازی، با بکار گیری هرس نمودن، بر پایه اصل کوتاه ترین توضیح می باشد. یکی دیگر از این تغییرات، تغییر درخت بگونه ای است که به جای آن که مشخص کند یک واحد متن می بایست در خلاصه قرار بگیرد یا خیر، احتمال انتخاب آن را برای قرار گرفتن در خلاصه تعیین می نماید. این الگوریتم برای تولید خلاصه از جنبه هایی چون طول، مکان و مجموعه هایی از ویژگی های زبانی استفاده می کند. شکل (۱) نمایش نمونه ای از یک درخت می باشد.



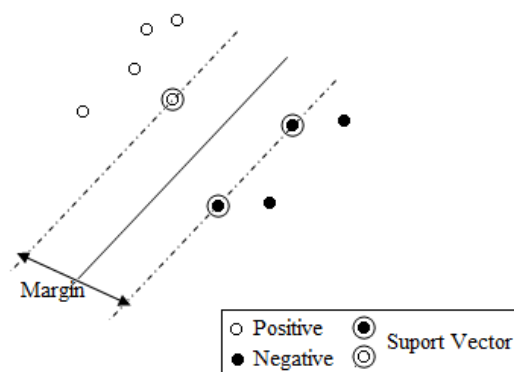
شکل (۱) - نمونه ای از یک درخت تصمیم<sup>۲</sup>

## • بردارهای پشتیبان<sup>۲</sup>:

در (ایران پور، ۱۳۸۶) پیشنهاد شده است که با در نظر گرفتن مجموعه ای از جنبه ها، چون مکان و طول برای جملات، یک ماشین بردار پشتیبان را بر روی یک پیکره خلاصه سازی آموزش داده و سپس به جملات بر اساس میزان فاصله آن ها از صفحه بردار پشتیبان امتیاز داده شود. در شکل (۲) نمونه ای از بردارهای پشتیبان نشان داده شده است.

<sup>۱</sup> . Decision tree

<sup>۲</sup> . Support Vector



شکل (۲) مثالی از بردارهای پشتیبان

## روش های ساختار کلامی

تئوری ساختار کلامی توسط تامسون<sup>۱</sup> و مان<sup>۲</sup> در سال ۱۹۸۸ ارائه گردید. ارتباط موجود بین دو بخش غیر هم پوشان متن به نام بخش پایه<sup>۳</sup> و بخش پیرو<sup>۴</sup> بیان می شود (مدرس خیابانی، ۱۳۸۵). تفاوت بین پایه و پیرو در این است که بخش پایه هدف اصلی نویسنده را بیان می دارد. مسئله دیگر آن است که بخش پایه و بخش پیرو از نظر معنایی و فهمی از یکدیگر مستقل می باشند. تئوری ساختار کلامی شامل ساختن درخت دودویی است که ارتباط انسجامی که در متن موجود است را بیان می دارد. ارتباط بین واحدهای بیانی تشخیص و برجسب گذاری می شود. دو هدف برای تحلیل ساختار کلامی وجود دارد. اول آن که بخش غیر مهم جمله ها حذف شود و دوم آن که ارتباط کلامی جمله های باقیمانده حفظ شود تا خلاصه منسجم ایجاد شود. سه سطح ارتباط کلامی پاراگراف، جمله و بخشی از جمله وجود دارد. ارتباط کلامی بین پاراگراف های همسایه به این صورت مشخص می شود که آیا آن ها یک موضوع را بر اساس ضریب باهم آیی اش دنبال می نمایند یا خیر. همچنین ارتباط کلامی بین جمله ها بر اساس استفاده از باهم آیی معنایی می باشد.

بر اساس تئوری ساختار زبانی یک متن منسجم می تواند به صورت یک درخت ساخت یافته در نظر گرفته شود که گره های میانی آن ارتباطات کلامی و برگ ها گزاره هایی می باشند که به وسیله بخش های موجود در متن توضیح داده می شوند. تجزیه کننده های تئوری ساختار زبانی وجود دارند که درخت تئوری ساختار زبانی مناسبی برای متن ایجاد می نمایند. مطابق با قاعده تئوری ساختار زبانی، تمامی واحدهای گزاره ای که در متن موجود می باشند، می بایست به وسیله ارتباط کلامی - بیانی<sup>۵</sup> به یکدیگر متصل باشند تا متن پیوسته باشد. اتصال واحدهای گزاره ای متن ساختار کلامی - بیانی را ایجاد می نماید. تئوری ساختار زبانی عموماً به وسیله درخت نشان داده می شود که هر ارتباطی در آن برای متصل کردن زیر درخت ها به کار می رود که خود می توانند درخت دیگری باشند.

در تئوری ساختار کلامی از یک سری نشانه ها همانند "اگرچه"، "بنابراین" و "در نتیجه" برای پیدا کردن ارتباط بین بخش های مختلف جمله استفاده می نمایند. ابتدا از یک تجزیه کننده استفاده می گردد تا مشخص شود چگونه زیردرخت ها به یکدیگر متصل شوند. ۲۵ ارتباط زبانی در تئوری ساختار زبانی وجود دارد که از آن جمله می توان به نتیجه و تصدیق

<sup>۱</sup> . Thompson

<sup>۲</sup> . Mann

<sup>۳</sup> . Nucleus

<sup>۴</sup> . Satellite

<sup>۵</sup> . Rhetorical

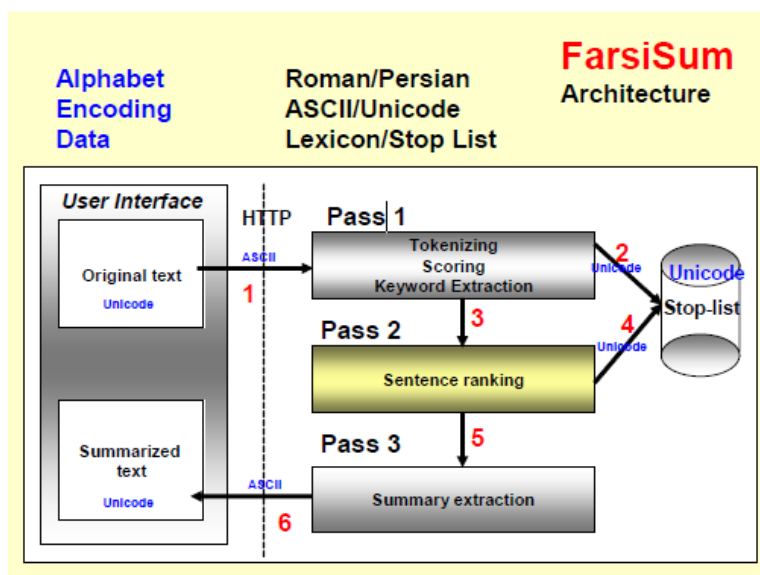
اشاره نمود. تمام روش هایی که بر اساس تئوری ساختار زبانی عمل می نمایند از روش سطحی بهتر هستند. این تئوری به میزان زیادی در محاسبات زبان شناختی استفاده می شود و دیگر کارهای انجام شده در حوزه پردازش زبان همانند تفکیک ضمائر و امتیاز دهی به سندها را تحت تأثیر قرار داده است (مشکی و آنالویی، ۱۳۸۷).

علاوه بر مشکل بودن روش هایی که بر اساس تحلیل ساختار متن عمل می نمایند، از مهم ترین معایب این روش ها می توان به پردازش های پرهزینه زبان شناختی اشاره نمود، این امر باعث می شود در مواردی نتیجه دلخواه ایجاد نشود. (Mani and Maybury, ۲۰۱۴)

لازم به ذکر است که کلیه روش های فوق الذکر را می توان به صورت ترکیبی جهت خلاصه سازی متون مورد استفاده قرار داد.

## مهم ترین پژوهش های انجام شده در خلاصه سازی در زبان فارسی

در این بخش به فعالیت های انجام شده در زبان فارسی پرداخته می شود. در (Hovy and Marcu, ۱۹۹۸) یک خلاصه ساز مبتنی بر وب به نام FarsiSum معرفی شده است که نسخه تغییر یافته یک خلاصه ساز سوئدی به نام SweSum برای پوشش زبان فارسی است. این سیستم در حال حاضر به صورت یک ابزار قابل استفاده برای کاربران می باشد. پروسه تولید خلاصه در این سیستم که در شکل (۳) نیز نشان داده شده است شامل سه مرحله ی ارسال درخواست خلاصه از طریق سرور HTTP، خلاصه سند در طی سه مرحله با استفاده از لیست توقف فارسی و بازگرداندن متن خلاصه شده از طریق HTTP به کاربر و پردازش متن خلاصه توسط مرورگر می باشد.



شکل (۳) FarsiSum

در مقاله (Barzilay et al, ۲۰۰۲) روشی برای ساخت سیستم خلاصه سازی خودکار متون فارسی مطرح کردند. کل روش به صورت نحوی/معنایی عمل می کند و ترکیبی از دو روش زنجیره های لغوی و خلاصه سازی مبتنی بر گراف است که از پنج معیار میزان شباهت جملات با یکدیگر، شباهت جملات با کلمات کلیدی کاربر، شباهت جملات با عنوان، تعداد

جملات مشابه هر جمله و وجود کلمات اشاره در جمله برای امتیازدهی به جملات استفاده نموده و جملات با بیشترین امتیاز را به عنوان خروجی سیستم انتخاب می کند. در واقع با محاسبه ارتباط بین جملات (گره های گراف) بر اساس زنجیره لغوی، شباهت بین جملات را بر اساس مفهوم آن ها در نظر گرفته اند. سیستم پیاده سازی شده براساس این روش با خلاصه های مرجعی که به صورت دستی تهیه شده اند، مورد مقایسه قرار گرفته اند. ارزیابی سیستم پیاده سازی شده به صورت دستی انجام گردیده است. به علت عدم دستیابی به منبعی شامل متون و خلاصه مرجع در زبان فارسی، خلاصه مرجع ۱۰ متن از روزنامه همشهری با نسبت فشردگی ۳۰ درصد توسط چند نفر از دانشجویان تولید شده است، که شباهت خلاصه مرجع با خلاصه سیستم در این متون به صورت متوسط ۵۶/۶۶٪ است. در (مشکی، ۱۳۸۸) یک روش مبتنی بر خوشه بندی برای خلاصه سازی چند سندی متون فارسی پیشنهاد گردیده است. در این روش، پس از پیش پردازش متن، در مرحله خلاصه سازی، ابتدا جمله ها خوشه بندی می شود و سپس به ازای هر خوشه جمله ای که بیشترین ارتباط با سایر جمله ها را دارد، گزینش می شود. در آخرین مرحله خلاصه سازی، جمله ها با توجه به ترتیب زمانی متن ها در خلاصه نهایی درج می گردند. نتایج پیاده سازی نشان می دهند که در بیشتر موارد خروجی سامانه خلاصه سازی پیشنهادی خلاصه قابل قبولی را تولید می کند. جهت ارزیابی نتایج این روش پیشنهادی نیز از شیوه قضاوت انسانی استفاده شده است. بدین صورت که برای خلاصه هر یک از ۱۰ مجموعه، ۳ داور انسانی رای خود را به صورت ۳۳٪ خوب، ۵۱٪ متوسط، ۱۶٪ ضعیف ارائه شده است. سیستم دیگری که سعی کرده معایب سیستم های مشابه موجود را برطرف نماید سیستم (شهاب، ۱۳۸۱) می باشد که ترکیبی از روش های مبتنی بر گراف و الگوریتم ژنتیک است که پس از وزن دهی جملات و تشکیل ماتریس شباهت برای جملات سند، گراف جهت داری بوجود می آورد. پس از آن سه فاکتور شباهت با عنوان، پیوستگی و قابلیت خوانایی برای جملات موجود در خلاصه محاسبه می شود و در مرحله بعد با استفاده از تابع برازندگی و الگوریتم ژنتیک، جملات خلاصه انتخاب می شوند. جهت ارزیابی این سیستم از ده متن در مقوله های مختلف علمی و خبری استفاده شده است. این متن ها توسط افراد مختلف با نسبت فشردگی مختلف خلاصه شده است. با مقایسه خلاصه سیستم و خلاصه های دستی مشاهده شده است که؛ به طور میانگین در ۶۰/۲۳ درصد موارد جملات خلاصه مانند جملات خلاصه دستی انتخاب شده اند. سیستم خلاصه ساز دیگری که در مقاله (مشکی، ۱۳۸۸) ارائه گردیده است با الهام از شیوه انسان در خلاصه سازی متن، روشی جدید برای خلاصه سازی متن براساس گزینش ارائه می دهد. مبنای روش پیشنهادی یافتن زنجیره ای از جمله ها می باشد که قوی ترین همبستگی را با هم داشته باشند. به این منظور متن به دو دید بیرونی و دید داخلی تقسیم بندی شده است. در دید بیرونی متن به پاراگراف ها تقسیم بندی می شود و در دید داخلی هر پاراگراف به مجموعه ای از جملات تقسیم می شود و اعضای داخلی هر پاراگراف را تشکیل می دهند. انتخاب پارامترها نقش مهمی را در مشخص کردن جمله های برجسته برای قرار گرفتن در خلاصه، بازی می کند. به همین دلیل در این روش ارائه شده از میزان نزدیکی جملات به اول یا آخر پاراگراف، تعداد کلمه های مشترک بین اجزای جمله ها، میزان فاصله پاراگراف ها از یکدیگر و تعداد کلمه های کلیدی مشترک بین اجزای جمله ها استفاده شده است. پس از انتخاب پارامتر میزان همبستگی بین جملات محاسبه شده و در نهایت زنجیره ای از جمله هایی را که قوی ترین ارتباط را با یکدیگر دارند پیدا می کند. در این زنجیره از هر پاراگراف فقط یک جمله انتخاب می شود و شرط انتخاب آن جمله قرار گرفتن در زنجیره ای می باشد که حاصل جمع همبستگی اجزای آن حداکثر باشد.

شروع زنجیره از پاراگراف اول می باشد و به دلیل این که الزاماً حق انتخاب یک جمله از هر پاراگراف وجود دارد، زنجیره نهایی از هر پاراگراف فقط یک عضو خواهد داشت. در واقع زنجیره حاصل که حاصل جمع همبستگی اجزای آن حداکثر است، همان خلاصه متن ورودی خواهد بود. به عبارت دیگر پس از پیدا کردن زنجیره مناسب به ازای هر عضو آن، جمله های متناظر انتخاب می شوند و خلاصه نهایی را تشکیل می دهند. برای ارزیابی نتایج حاصل از روش ذهنی استفاده شده است. در این نوع ارزیابی از افراد خبره به عنوان داور برای قضاوت در مورد خلاصه تولید شده استفاده می شود. برای این منظور ۴ داور از افراد خبره انتخاب شده و نتایج خلاصه سازی به همراه اصل متون را بررسی کرده اند و قضاوت خود را در رابطه با نتایج بدست آمده در قالب یکی از کلمات ضعیف، قابل قبول و خوب بیان کرده اند. که ۴۰٪ رای خوب ۴۵٪ قابل قبول و ۱۵٪ ضعیف داده شده است.

### نتیجه گیری

به طور کلی روش های خلاصه سازی در چهار دسته قابل تقسیم بندی هستند که شامل روش های آماری سطحی، درک متن، مبتنی بر یادگیری نظارت شده و ساختار کلامی هستند. روش های آماری سطحی تنها از ویژگی های آماری از قبیل فراوانی واژه ها، مکان واحد های متنی، عبارات توضیحی، امضاء سرفصل و با هم آیی برای تولید خلاصه استفاده می کنند (Barzilay et al, ۲۰۰۲). مزیت اصلی این گونه روش ها این است که بسیار ساده اند و بر روی هر نوع متن ورودی می توان آن ها را اعمال کرد و عیب مهم این دسته از روش ها عدم در نظر گرفتن ارتباط بین کلمات است و به همین دلیل عدم انسجام و پیوستگی در خلاصه های تولیدی به وسیله ی این روش ها به وضوح قابل مشاهده است. دسته دوم از روش های خلاصه سازی، روش های درک متن می باشند. اینگونه روشها علاوه بر نیاز به منبع دانش نیازمند تفسیر پیچیده متن می باشند ولی عملکرد بهتری را نسبت به روش های آماری ارائه می دهند. در این روش ها از ارتباط میان موجودیت ها شامل شباهت، هم معنی، ارتباط نحوی و ارتباط معنایی، ارتباط لغوی و هم مکانی استفاده می کنند. دسته سوم روش های مبتنی بر یادگیری نظارتی بودند که گونه ای از یادگیری است که در آن به عامل یادگیرنده تعدادی جنبه و مقدار متناظر برای خصوصیتی یاد گرفته شود داده شده و آموزش لازم جهت عمل نمودن در موارد مشابه داده می شود و دسته ی آخر روش های ساختار کلامی بودند که برای تشخیص بخش های مهم متن استفاده می شدند. توجه به این روش باعث انسجام خلاصه های تولیدی می شود.

با توجه به موارد ذکر شده در این پژوهش، توجه به ترکیب روش های مطرح شده برای تولید خلاصه های باکیفیت تر بسیار اهمیت دارد. توجه به یکی از این روشها به تنهایی باعث می شود که خلاصه تولیدی از لحاظ پیوستگی و انسجام و یا انتخاب جملات دچار مشکل گردد. به طور مثال در روش های آماری تنها تکیه بر مباحث آماری شده است و در صورتی که با روش های درک متن و ارتباط کلامی ترکیب گردد به نتایج قابل قبول تری می رسد. همچنین می توان یک ضریبی را تعیین نمود و مقدار تاثیر هر کدام از روش ها را با استفاده از آن ضریب تنظیم نمود و نقطه طلایی سیستم را که در آن حالت خلاصه های بهتری تولید می گردد را تعیین نمود. کارهایی که در زبان فارسی انجام گرفته بیشتر به روش های آماری و مبتنی بر درک متن (شامل روشهای خوشه بندی و مبتنی بر گراف) متکی است حال آنکه از طریق ترکیب با روش های روابط کلامی می تواند به نتایج بهتری رسید.

## منابع و مراجع

- شهاب، امیر شهاب، چکیده سازی متون فارسی، دومین کنفرانس علوم شناختی، تهران، ص ۵۶، ۱۳۸۱.
- کریمی، زهره، شمس فرد، مهرنوش، سیستم خلاصه سازی خودکار متون فارسی، دوازدهمین کنفرانس بین المللی انجمن کامپیوتر ایران، تهران، دانشگاه شهید بهشتی، ۱۳۸۵.
- مشکی محسن، آنالوئی، مرتضی، خلاصه سازی چند سندی متون فارسی با استفاده از یک روش مبتنی بر خوشه بندی، اولین کنفرانس ملی مهندسی نرم افزار، دانشگاه آزاد رودهن، ۱۳۸۸.
- اخوان، تارا، شمس فرد، مهرنوش. عرفانی جورابچی، مونا، PARSUMIST خلاصه ساز تک سندی و چند سندی متون فارسی، چهاردهمین کنفرانس ملی سالانه انجمن کامپیوتر ایران، تهران، دانشگاه صنعتی امیر کبیر، ۱۳۸۷.
- بهره پور، مجید، مهدی پور، الهام، کامل، آزاده، امیری، ملیحه، طهماسبی، آیدا، اکبرزاده توتونچی، محمدرضا، سیستم خلاصه سازی خودکار متن های فارسی، چهاردهمین کنفرانس ملی سالانه انجمن کامپیوتر ایران، تهران، دانشگاه صنعتی امیر کبیر، ۱۳۸۷.
- ستوده، حمیدرضا، اکبرزاده توتونچی، محمدرضا، تشنه لب، محمد، خلاصه سازی متن بر اساس گزینش با استفاده از رویکرد انسان شناختی، هجدهمین کنفرانس مهندسی برق ایران، اصفهان، دانشگاه اصفهان، صفحه ۲۱-۲۳، ۱۳۸۹.
- بازقندی، مهدی، تدین تبریزی، قمرناز، وفایی جان، مجید، بازقندی، علی، خلاصه سازی گزینشی متون فارسی مبتنی بر خوشه بندی PSO، اولین کنفرانس بین المللی پردازش خط و زبان فارسی، ایران، دانشگاه سمنان، ۱۳۹۱.
- هنریشه، محمد علی، طراحی و پیاده سازی یک سیستم خلاصه ساز متون فارسی، پایان نامه دوره کارشناسی ارشد، دانشگاه صنعتی شریف، دانشکده مهندسی کامپیوتر، ۱۳۸۶.
- ارژنگ، غلامرضا، زبان و ادب فارسی، تهران: نشر قطره، ۱۳۸۳.
- ریاحی، نوشین، غزالی، فاطمه، محمد علی، بهبود کارایی سیستم خلاصه ساز متون فارسی با استفاده از الگوریتم هرس در شبکه های عصبی، اولین کنفرانس بین المللی پردازش خط و زبان فارسی، ایران، دانشگاه سمنان، ۱۳۹۳.
- حافظی، محمد مهدی، ثامتی، حسین، منصوری، نیلوفر، منتظری، نیلوفر، بحرانی، محمد، موثق، حامد، ارائه یک مدل دستوری برای بهبود دقت سیستم های بازشناسی گفتار پیوسته فارسی، دومین کارگاه پژوهشی زبان فارسی و رایانه، تهران، دانشگاه تهران، ۱۳۸۵.
- حسامی فرد، رضا، قاسم ثانی، غلامرضا، طراحی و پیاده سازی یک الگوریتم ریشه یابی برای زبان فارسی، یازدهمین کنفرانس بین المللی انجمن کامپیوتر ایران، تهران، ص ۵۱۵-۵۱۹، ۱۳۸۴.
- عربی نرئی، سمیه، وحیدی اصل، مجتبی، مینایی بیدگلی، بهروز، استخراج کلمات کلیدی جهت طبقه بندی متون فارسی، اولین کنفرانس داده کاوی ایران، تهران، دانشگاه صنعتی امیر کبیر، ۱۳۸۶.
- ایران پور مبارکه، مهدی، بررسی مشکلات تعیین حدود جمله و کلمه، سمینار کارشناسی ارشد، دانشگاه علم و صنعت، ۱۳۸۶.
- مدرس خیابانی، شهرام، قیومی، مسعود، نقش پیکره های زبانی در با هم آیی: رویکردی مقایسه ای، دومین کارگاه پژوهشی زبان فارسی و رایانه، تهران، ص ۵۵ تا ۵۶، ۱۳۸۵.
- مشکی، محسن، آنالوئی، مرتضی، یک روش آماری مبتنی بر پیکره برای جداسازی واژه های به هم چسبیده، دومین کنگره مشترک سیستم های فازی و هوشمند ایران، ۱۳۸۷.
- مشکی، محسن، خلاصه سازی گزینشی چند سندی متون فارسی، پایان نامه دوره کارشناسی ارشد، دانشگاه علم و صنعت ایران، دانشکده مهندسی کامپیوتر، ۱۳۸۸.



فتوحی، رضا، ارزیابی کارایی الگوریتم های رتبه بندی صفحه برای استخراج صفحات وب، هشتمین سمپوزیوم مهندسی کامپیوتر و توسعه پایدار با محوریت شبکه های کامپیوتری، مدل سازی و امنیت سیستم ها، موسسه آموزش عالی خاوران، مشهد ۱۳۹۲.

McKeown, k., Barzilay, R., Chen, J., Elson, D., Evans, D., Klavans, J., Nenkova, A., Schiffman, B., and Sigelman, S. (۲۰۰۳). "Columbia'S News Blaster: New Features and Future Directions". In North American Chapter of the Association for Computational Linguistics on Human Language Technology, Pages ۱۵-۱۶.

Buyukkokten, O., Gacia – Molina, H., and Paepke A. (۲۰۰۱). "seeing the Whole in Parts: Text Summarization for web Browsing on Handheld Devices". In ۱۰th international www Conference. Hong Kony, China.

Elhadad, N., Kan, M., Klavans, J., and McKeown, K.( ۲۰۰۵) "Customization in a unified framework for summarizing medical literature". In Journal Artificial Intelligence in Medicine, Volume ۳۳, Pages ۱۷۹-۱۹۸.

Mazdak, N., Hassel, M. (۲۰۰۴). "FarsiSum-a persian text summarizer", Master thesis, Department of linguistics, Stockholm University .

Dalianis, H. (۲۰۰۰). "SweSum - A Text Summarizer for Swedish, Technical report ."In Interaction and Presentation Laboratory, TRITA-NA-P۰۰۱۵, IPLab-۱۷۴

Barzilay, R., Elhadad, N., and McKeown, K. (۲۰۰۲). "Inferring Strategies for Sentence Ordering in Multidocument News Summarization". In Journal of Artificial Intelligence Research, Pages ۵۵-۳۵

Chuang, W., Yang, J. (۲۰۰۰). "Extracting Sentence Segments for Text Summarization: A machin Learning Approach". In ۲۳rd annual international Conference ACM SIGIR, Athens, Greece.

Porter, M. F. (۱۹۸۰). "Readings in information retrieval". Morgan Kaufmann Publishers Inc. San Francisco, Pages ۳۱۶-۳۱۳

Mani, I., and Maybury, M. ( ۲۰۱۴). Advances in Automatic Text Summarization .Cambridge, MA: The MIT Press .

Hovy, E., and Marcu, D. (۱۹۹۸). "Automated Text Summarization". Tutorial at OLING/ACL'۹۸

R. Ferreira, L. de Souza Cabral, F. Freitas, R. D. Lins, G. de França Silva, S. J. Simske, et al. ۲۰۱۴, "A multi-document summarization system based on statistics and linguistic treatment," Expert Systems with Applications, vol. ۴۱, pp. ۵۷۸۷-۵۷۸۰.