

خلاصه‌سازی متن بر اساس گزینش با استفاده از رویکرد انسان‌شناختی

حمیدرضا ستوده¹، محمدرضا اکبرزاده توتونچی² و محمد تشنه‌لب³

¹گروه کامپیوتر، دانشگاه آزاد اسلامی واحد علوم و تحقیقات تهران، H.setoodeh@gmail.com

²گروه برق، دانشکده مهندسی، دانشگاه فردوسی مشهد، Akbarzadeh@ieee.org

³دانشکده برق، دانشگاه خواجه نصیرالدین طوسی، Teshnehlab@eed.kntu.ac.ir

چکیده - در این مقاله با الهام از شیوه انسان در خلاصه‌سازی متن، روشی جدید برای خلاصه‌سازی متن بر اساس گزینش ارائه شده است. در این روش ابتدا پارگراف‌ها و جمله‌های تشکیل دهنده متن مشخص شده و کلمه‌های غیرضروری و حروف اضافه حذف می‌شوند. سپس دسته‌هایی از جمله‌ها تشکیل شده و بر اساس یکسری از ویژگی‌ها ارتباط‌هایی بین جمله‌های داخلی و خارجی تمام دسته‌ها شکل می‌گیرد که منجر به ایجاد دنباله‌ای از جمله‌های برجسته به عنوان خلاصه می‌گردد.

کلید واژه - خلاصه‌سازی متن، گزینش

1- مقدمه

خلاصه‌سازی متن فرایندی امکان پذیر است زیرا به طور معمول در یک متن افزونگی رخ می‌دهد یعنی توضیحاتی اضافه وجود دارد که می‌توان آنها را خلاصه یا حذف کرد. همچنین اطلاعات مهم که در برگیرنده معنی و مفهوم اصلی متن می‌باشند نیز به طور یکنواخت در سطح متن پخش نشده‌اند.

تعریف یکسانی برای مهم و برجسته بودن و همچنین افزونگی در متن وجود ندارد زیرا انسان‌ها دارای اطلاعات و علائق مختلفی می‌باشند، و ممکن است موضوعی برای شخصی بدیهی باشد در عین حال همان مطلب برای فرد دیگری از اهمیت خاصی برخوردار باشد. بنابراین یک قضاوت ثابت و پایدار در مورد کیفیت یک خلاصه بسیار سخت است و این واقعیت ارزیابی خلاصه‌سازی را با مشکل مواجه ساخته است [2].

تنوع روش‌های خلاصه‌سازی به نحوه معرفی اجزای برجسته مربوط می‌باشد. روش‌های خلاصه‌سازی متن از همین دید به دو رویکرد خلاصه‌سازی بر اساس گزینش و خلاصه‌سازی بر اساس مفهوم تقسیم‌بندی می‌شوند.

در رویکرد خلاصه‌سازی بر اساس گزینش معمولاً از اطلاعات آماری مانند تعداد تکرار کلمه‌ها، کلمه‌های به کار رفته در عنوان، مکان جمله‌ها و طول جمله‌ها در متن برای استخراج جمله‌ها و اجزای برجسته استفاده می‌شود [3].

در روش خلاصه‌سازی بر اساس گزینش معمولاً از

با شکل‌گیری انقلاب اطلاعاتی، روزانه هزاران سند متنی تولید می‌شوند که این اسناد و مدارک به عنوان رسانه اصلی در کسب و کار و انتشار اطلاعات علمی مطرح می‌باشند. این اسناد حجم بالایی دارند و امکان تجزیه و تحلیل همه‌ی این سندها برای انسان وجود ندارد. بنابراین باید به دنبال روشی باشیم که انسان را در بهره‌برداری و استفاده بهتر از این گونه اطلاعات کمک کند. به منظور استفاده موثر و کامل از این اطلاعات، استخراج اجزای اصلی یا معنای کلی این اسناد در حجمی کمتر نسبت به اطلاعات اصلی امری ضروری می‌باشد.

با توجه به مطالب فوق وجود سیستمی خودکار که توانایی خلاصه‌سازی متن را داشته باشد و انسان را در تجزیه و تحلیل این اسناد یاری کند بسیار مفید می‌باشد.

خلاصه‌سازی متن فرایند شناسایی و استخراج برجسته‌ترین اجزای یک سند یا مجموعه‌ای از اسناد و گردآوری آنها در حجمی کمتر می‌باشد که این اجزا بیشترین تطابق و تعلق به موضوع و مفهوم مطرح شده در سند را داشته باشند [1].

به منظور تولید خلاصه باید به شناسایی اجزای برجسته در اطلاعات موجود در سند یا متن پردازیم، اطلاعات غیر مرتبط را به حداقل برسانیم، جزئیات را کاهش دهیم و آنها را در یک گزارش منسجم گردآوری کنیم.

تکنیک‌های سطحی متن استفاده می‌شود. یعنی ویژگی‌های آماری در هر قسمت از متن مشخص شده سپس بر اساس آن ویژگی‌ها درجه برجستگی هر جمله از متن مشخص می‌شود و در نهایت جمله‌هایی که بالاترین درجه برجستگی را کسب کنند، در خلاصه‌ی نهایی قرار می‌گیرند.

در مقایسه با رویکرد قبل در خلاصه‌سازی بر اساس مفهوم، علاوه بر اطلاعات آماری، از تجربه و تحلیل لغوی و معنایی متن برای ایجاد یک خلاصه با درجه بالای پیوستگی استفاده می‌شود [4].

خلاصه‌سازی بر اساس مفهوم، به خاطر دشواری‌های مربوط به پردازش زبان طبیعی همچنین نبود ابزارهای تجزیه و تحلیل برای همه‌ی زبان‌های طبیعی با مشکلاتی مواجه است و معمولاً به طور کامل امکان‌پذیر نمی‌باشد.

خلاصه‌سازی ممکن است بر روی یک سند یا چندین سند انجام شود که تمرکز در این مقاله بر روی خلاصه‌سازی بر اساس یک سند می‌باشد.

در این مقاله یک روش مبتنی بر گزینش برای خلاصه‌سازی متن ارائه شده است. در این روش پس از انجام کارهای مقدماتی بر روی متن، پاراگراف‌ها و جمله‌های تشکیل دهنده‌ی هر پاراگراف در متن مشخص می‌شوند. سپس دسته‌هایی از جمله‌ها تشکیل شده و بر اساس چند ویژگی، ارتباط‌هایی بین جمله‌های درونی و بیرونی دسته‌ها شکل می‌گیرد و منجر به ایجاد دنباله‌ای از جمله‌های برجسته به عنوان خلاصه می‌گردد.

2- مروری بر کارهای گذشته

مدل فضای برداری یکی از روش‌های معمول برای نشان دادن یک متن در فضای برداری می‌باشد. در این مدل متن به صورت یک بردار از میزان تکرار اجزا یا معکوس تکرار اجزا نشان داده می‌شود. سپس جمله‌هایی که اجزایشان بیشترین تعداد تکرار را دارند انتخاب شده و در خلاصه نهایی آورده می‌شوند. در [5] مدل تکرار اجزا برای خلاصه‌سازی متن در نظر گرفته شده است.

روش‌هایی پیشنهاد شده در [6] و [7] بر این اساس است که اگر وزن هر کدام از اجزا و کلمه‌ها بدون در نظر گرفتن میزان برجستگی و اهمیت جمله‌های در برگیرنده آن‌ها محاسبه شود به خلاصه مناسبی منجر خواهد شد. بنابراین بهتر است وزن هر کدام از اجزا را با در نظر گرفتن میزان مهم بودن جمله‌های آن‌ها محاسبه شود.

بعضی از روش‌های خلاصه‌سازی براساس گزینش، بر روی کلمات کلیدی تکیه دارند که اصطلاحاً به آن «جان کلام» گفته می‌شود.

«جان کلام» در حقیقت به موضوع اصلی مورد نظر نویسنده و یا مفهوم اصلی که خواننده از متن درک می‌کند، گفته می‌شود و ممکن است برای یافتن آن در متن از یک پایگاه داده خارجی نیز استفاده شود [8].

اگر از روش‌های آماری ساده استفاده کنیم «جان کلام» مهمترین قسمت از یک متن را تشکیل می‌دهد که به وسیله یک جمله بیان می‌شود. در حقیقت می‌توان گفت «جان کلام» یک راهنما برای یافتن سایر جمله‌های مناسب برای خلاصه می‌باشد. سایر جمله‌ها به خلاصه اضافه می‌شوند در صورتی که میزان همبستگی لازم با «جان کلام» را داشته باشند [9].

در [10] یک روش خلاصه‌سازی متن با استفاده از یک سیستم استنتاج فازی مطرح شده است که در این روش از منطق فازی برای اندازه‌گیری درجه برجستگی و همچنین مشخص کردن جمله‌های مهم برای ایجاد خلاصه، استفاده شده است.

در [11] روشی بهینه برای خلاصه‌سازی متن با استفاده از سیستم استنتاج فازی پیشنهاد شده است که توابع عضویتی که در سیستم فازی آن که برای انتخاب جمله‌های متن به کار رفته است، توسط الگوریتم ژنتیک بهینه شده و جواب‌های قابل قبول‌تری نسبت به [10] در بر داشته است.

سایر روش‌هایی که در خلاصه‌سازی بر اساس گزینش به کار رفته است به نوعی بهینه‌سازی روش‌های فوق می‌باشند.

3- روش پیشنهادی

معمولاً در روش‌های خلاصه‌سازی مبتنی بر گزینش خلاصه بر اساس یافتن مهمترین اجزای متن و گردآوری آن‌ها در کنار یکدیگر شکل می‌گیرد.

اما به دلیل آنکه معمولاً سندهایی که به عنوان متن ورودی سیستم‌های خلاصه‌ساز مطرح می‌باشند، به نحوی توسط انسان تولید شده‌اند، می‌توان گفت که یک زنجیره از جمله‌ها در متن وجود دارد که مفهوم و مقصود اصلی نویسنده را در بر می‌گیرند و سایر جمله‌ها در سایه این زنجیره قرار می‌گیرند. بنابراین برای ایجاد خلاصه فقط یافتن مهمترین اجزای متن کافی نمی‌باشد و باید به پیوستگی جمله‌ها بیشتر توجه داشت.

بنابراین در این مقاله با الهام از شیوه انسان در تفکر روشی پیشنهاد شده است که برای ایجاد خلاصه به جای یافتن

مهمترین جمله‌ها، زنجیره و دنباله‌ای از جمله‌ها در کل متن انتخاب شوند که با هم همبستگی و ارتباط بیشتری داشته باشند.

حال ممکن است هر یک از این جمله‌های زنجیره به تنهایی در کل متن مهم نباشند ولی با قرار گرفتن در زنجیره‌ای از جمله‌ها پازلی را تکمیل کنند که این پازل در کل نشانگر مفهوم و منظور متن ابتدایی باشد. همان‌طور که در یک پازل ممکن است هر کدام از قطعه‌ها به تنهایی مهم نباشند ولی وقتی در پازل قرار می‌گیرند تکمیل کننده سایر قطعات چیده شده خواهند بود.

بنابراین مبنای روش پیشنهادی ما یافتن زنجیره‌ای از جمله‌ها می‌باشد که قوی‌ترین همبستگی را با هم داشته باشند. در روش ارائه شده متن را به دو دید بیرونی و دید داخلی تقسیم بندی کرده‌ایم. در دید بیرونی متن به پاراگراف‌ها تقسیم بندی می‌شود و در دید داخلی هر پاراگراف به مجموعه‌ای از جمله‌ها تقسیم می‌شود و اعضای داخلی هر پاراگراف را تشکیل می‌دهند.

3-1- انتخاب پارامترها

انتخاب پارامترها نقش مهمی را در مشخص کردن جمله‌های برجسته برای قرار گرفتن در خلاصه، بازی می‌کند. به همین دلیل در این مقاله از میزان نزدیکی جمله‌ها به اول یا آخر پاراگراف، تعداد کلمه‌های مشترک بین اجزای جمله‌ها، میزان فاصله پاراگراف‌ها از یکدیگر و تعداد کلمه‌های کلیدی مشترک بین اجزای جمله‌ها استفاده شده است. این 4 پارامتر بصورت زیر محاسبه می‌شوند:

1. میزان نزدیکی جمله‌ها به اول یا آخر پاراگراف

معمولاً جمله‌های اول و آخر هر پاراگراف شامل نکات مهمتری نسبت به سایر جمله‌ها می‌باشند. بنابراین اگر اجزای پاراگراف‌ها جزء جمله‌های اول و آخر هر پاراگراف باشند داری ارتباط قوی‌تری نسبت به سایر جمله‌ها می‌باشند.

ارتباط جمله‌های هر پاراگراف با سایر پاراگراف‌ها به صورت زیر محاسبه می‌شود:

$$A(P_{ij} + P_{mn}) = \bar{X}P_{ij} \times \bar{X}P_{mn} \quad i \neq m \quad (1)$$

که در آن $\bar{X}P_{ij}$ و $\bar{X}P_{mn}$ محل نسبی قرار گرفتن جمله‌های j و i در پاراگراف‌های m و i می‌باشند و

این مقدار برای جمله j از پاراگراف i بصورت زیر محاسبه می‌شود:

$$\bar{X}P_{ij} = \frac{|CP_i - XP_{ij}|}{CP_i - 1} \quad (2)$$

$$XP_{ij} \in [1, 2, \dots, n_i] \quad (3)$$

که CP_i در رابطه 2 مرکز پاراگراف i می‌باشد و بصورت زیر تعریف می‌شود:

$$CP_i = \frac{n_i + 1}{2} \quad (4)$$

و XP_{ij} در رابطه 2 و 3 مکان قرار گرفتن جمله j در پاراگراف i و n_i در رابطه 3 و 4 برابر با تعداد جمله‌های پاراگراف i می‌باشد.

2. تعداد کلمه‌های مشترک بین اجزای جمله‌ها

اجزایی از پاراگراف‌ها که کلمه‌های مشترک بیشتری داشته باشند دارای همبستگی و ارتباط قوی‌تری نسبت به سایر اجزا می‌باشند. برای بدست آوردن میزان این ارتباط جمله‌های هر پاراگراف به کلمه‌ها تقسیم می‌شوند و پس از انجام اعمال مقدماتی و حذف کلمه‌های اضافه، بصورت دوبعدی تعداد کلمه‌های مشترک بین جمله‌های هر پاراگراف را بدست می‌آوریم.

برای اینکه جمله‌های بلند شانس بیشتری نسبت به جمله‌های کوچک نداشته باشند تعداد کلمه‌های مشترک دو جمله را بر میانگین کل تعداد کلمه‌های آن جمله‌ها تقسیم می‌کنیم. این پارامتر بصورت زیر تعریف می‌شود:

$$F(P_{ij}, P_{mn}) = \frac{N_{jn}(w)}{n_i + n_m} \quad i \neq m \quad (5)$$

که در آن $F(P_{ij}, P_{mn})$ میزان ارتباط بین جمله‌های j و n از لحاظ کلمه‌های مشترک و $N_{jn}(w)$ تعداد کلمه‌های مشترک جمله‌های j و n از پاراگراف‌های i و m و همچنین n_i و n_m برابر با تعداد جمله‌های این پاراگراف‌ها می‌باشند.

3. تعداد کلمه‌های کلیدی مشترک بین جمله‌ها

جمله‌های از پاراگراف‌ها که تعداد کلمه‌های کلیدی مشترک بیشتری دارند و یا این که شامل کلمه‌های به کار رفته در عنوان متن می‌باشند دارای اهمیت و همبستگی

بیشتری نسبت به سایر جمله‌ها می‌باشند.

پس از مشخص کردن کلمه‌های کلیدی متن و کلمه‌های موجود در عنوان، میزان تشابه دبدو بین جمله‌ها را بدست می‌آوریم. کلمه‌های کلیدی کلمه‌هایی هستند که وزن تکرار آنها در متن از سایرین بیشتر باشد. برای بدست آوردن میزان تشابه د، جمله از لحاظ کلمه‌های کلیدی به صورت زیر عمل می‌شود:

$$\bar{F}(P_{ij}, P_{mn}) = \lambda \times \frac{N_{jm}(\bar{w})}{n_i + n_m} \quad i \neq m \quad (6)$$

که در آن $\bar{F}(P_{ij}, P_{mn})$ میزان ارتباط بین جمله‌های j و n از لحاظ کلمه‌های کلیدی و $N_{jm}(\bar{w})$ تعداد کلمه‌های مشترک کلیدی یا عنوان در جمله‌های j و n و همچنین n_i و n_m برابر با تعداد جمله‌های پاراگراف‌های i و m می‌باشند. λ ضریبی است که برای این پارامتر بدلیل میزان مهمی آن در نظر گرفته‌ایم و مقدار آن در سیستم $1/5$ منظور شده است.

4. میزان فاصله پاراگراف‌ها از یکدیگر

در یک متن پاراگراف‌های نزدیک به یکدیگر دارای ارتباط قوی‌تری نسبت به سایر پاراگراف‌ها می‌باشند و به همین ترتیب اجزای تشکیل دهنده این پاراگراف نیز دارای ارتباط قوی‌تری نسبت به یکدیگر می‌باشند. که این ارتباط به صورت زیر محاسبه می‌شود:

$$D(P_i, P_m) = \frac{1}{|i - m|} \quad i \neq m \quad (7)$$

که در آن $D(P_i, P_m)$ میزان ارتباط اجزای پاراگراف‌های i و m از لحاظ میزان فاصله نسبت به یکدیگر می‌باشد.

3-2- محاسبه همبستگی جمله‌ها

پس از محاسبه این 4 پارامتر میزان همبستگی بین جمله‌ها بصورت زیر معین می‌شود:

$$R(P_{ij}, P_{mn}) = A(P_{ij}, P_{mn}) + F(P_{ij}, P_{mn}) + \bar{F}(P_{ij}, P_{mn}) + D(P_i, P_m) \quad (8)$$

که در آن $R(P_{ij}, P_{mn})$ میزان همبستگی و ارتباط بین جمله‌های j و n از پاراگراف‌های i و m می‌باشد.

در رابطه 8 شرط $i \neq m$ برقرار می‌باشد به‌دلیل این‌که میزان همبستگی باید بین اجزای دو پاراگراف مختلف محاسبه شود.

3-3- ایجاد زنجیره نهایی

پس از مشخص شدن ارتباط بین تمام جمله‌های پاراگراف‌های مختلف، زنجیره‌ای از جمله‌هایی را که قوی‌ترین ارتباط را با یکدیگر دارند پیدا می‌کنیم. در این زنجیره از هر پاراگراف فقط یک جمله انتخاب می‌شود و شرط انتخاب آن جمله قرار گرفتن در زنجیره‌ای می‌باشد که حاصل جمع همبستگی اجزای آن حداکثر باشد.

بنابراین جمله j از پاراگراف i به شرطی در زنجیره نهایی قرار می‌گیرد که زنجیره وجود داشته باشد که حاصل رابطه زیر حداکثر شود:

$$S = \dots + R(P_{mn}, P_{ij}) + R(P_{ij}, P_{xy}) + \dots \quad (9)$$

شروع زنجیره از پاراگراف اول می‌باشد و به دلیل این‌که فقط حق انتخاب یک جمله از هر پاراگراف را داریم، زنجیره نهایی از هر پاراگراف فقط یک عضو خواهد داشت.

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، جمله‌های متناظر انتخاب می‌شوند و خلاصه نهایی را تشکیل می‌دهند.

پیاده‌سازی مراحل ذکر شده بوسیله زبان برنامه‌نویسی Visual C# انجام شده است.

4- نتیجه‌گیری

با توجه به این‌که می‌توان گفت متن از تفکر انسانی ناشی شده است بنابراین با الهام از شیوه انسان در خلاصه‌سازی بهتر می‌توان به بحث در مورد خلاصه‌سازی متن و شیوه برخورد با آن پرداخت. به همین منظور در این مقاله روشی ارائه شده است که خلاصه‌سازی را بر مبنای یافتن زنجیره‌ای از جمله‌ها انجام می‌دهد که به نظر می‌رسد این زنجیره همان دنباله از افکار نویسنده‌ای بوده است که در نوشتن متن نقشی را ایفا کرده است.

انتخاب ویژگی‌ها نقش مهمی را در مشخص کردن جمله‌های برجسته برای قرار گرفتن در خلاصه، بازی می‌کند. به همین دلیل در این مقاله از میزان نزدیکی جمله‌ها به اول یا آخر پاراگراف، تعداد کلمه‌های مشترک بین اجزای جمله‌ها، میزان فاصله پاراگراف‌ها از یکدیگر و تعداد کلمه‌های کلیدی مشترک بین اجزای جمله‌ها استفاده شده است. برای بهینه کردن روش پیشنهادی می‌توان از تعداد بیشتری از ویژگی‌ها بین جمله‌ها بهره برد. هر چه تعداد پارامترها بیشتر باشد ممکن است مقدار

Conf. on Research and Development in Information Retrieval. Greece, pp. ۱۲۵-۱۵۹, ۲۰۰۰.

- [A] P. Filho and T. Pardo, "Summarizing Scientific Texts: Experiments with Extractive Summarizers," Proc. Int. Conf. on Intelligent Systems Design and Applications. Brazil, pp. ۵۲۰-۵۲۴, ۲۰۰۷.
- [۹] T. A. S. Pardo, L. H. M. Rino and M. G. V. Nunes, "GistSumm: A Summarization Tool Based on a New Extractive Method," Proc. Int. Conf. on Artificial Intelligence. Portugal, pp. ۱-۱۰, ۲۰۰۳.
- [۱۰] A. Kiani and M.R. Akbarzadeh T., "Intelligent Extractive Text Summarization Using Fuzzy Inference Systems," Proc. Int. Conf. on Intelligent Engineering. Canada, pp. ۲۷۰-۲۷۸, ۲۰۰۶.
- [۱۱] A. Kiani and M.R. Akbarzadeh T., "Automatic Text Summarization Using Hybrid Fuzzy GA-GP," Proc. Int. Conf. on Fuzzy Systems. Canada, pp. ۲۷۰-۲۷۸, ۲۰۰۶.

دقیق‌تری برای همبستگی جمله‌ها حاصل شود و خلاصه نهایی بهتر بتواند کل متن را پوشش دهد.

روش پیشنهادی بر روی مجموعه‌ای از متون خبری موجود در اینترنت پیاده‌سازی شد. این متون به گونه‌ای انتخاب شدند که هر یک شامل پاراگراف‌های مختلفی باشند. و هر کدام از این متون به طور جداگانه به عنوان ورودی به سیستم خلاصه‌ساز داده شدند.

برای ارزیابی نتایج حاصل از روش ذهنی (subjective) استفاده شد. در این نوع ارزیابی از افراد خبره به عنوان داور برای قضاوت در مورد خلاصه تولید شده استفاده می‌شود.

برای این منظور 4 داور از افراد خبره انتخاب شدند و نتایج خلاصه‌سازی به همراه اصل متون را بررسی کردند و قضاوت خود را در مورد نتایج بدست آمده در قالب یکی از کلمات ضعیف، قابل قبول و خوب بیان کردند. نتایج حاصل در جدول 1 آورده شده است.

جدول 1: نتایج ارزیابی خلاصه‌سازی

ضعیف	قابل قبول	خوب
٪15	٪45	٪40

هر چند در اینجا از روش ذهنی برای ارزیابی استفاده شده است و نمی‌توان به طور قطعی در مورد مقایسه با سایر روش‌ها سخن گفت اما در مقایسه با سایر روشها مانند روش برداری [5] و روش «جان کلام» [8]، خلاصه تولید شده در این روش بدلیل اینکه در انتخاب جمله‌ها اهمیت بیشتری به پیوستگی در مقابل برجستگی جمله‌ها داده شده، دارای پیوستگی بیشتری می‌باشد.

مراجع

- [1] I. Mani and M. Maybury, *Advanced in automatic summarization*. John Benjamins Publishing Company, pp. ۱۲۹-۱۶۵, ۲۰۰۱.
- [۲] J. Kupiec, J. Pedenon and F. Cheb, "A Trainable Document Summarizer," Proc. Int. Conf. on Research and Development in Information Retrieval. Washington, pp. ۶۸-۷۳, ۱۹۹۵.
- [۳] J. Goldstein, M. Kantrowitz and V. Mittal, "Summarizing text documents: Sentence selection and evaluation metrics," Proc. Int. Conf. on Research and Development in Information Retrieval. California, pp. ۱۲۱-۱۲۸, ۱۹۹۹.
- [۴] L. Dey, A. C. Rastogi and S. Kumar, "Generating Concept Ontologies through Text Mining," Proc. Int. Conf. on Web Intelligence (WI'۰۶), pp. ۲۳-۲۳, ۲۰۰۶.
- [۵] K. B. Khoo and M. Ishizuka, "Topic Extraction from News Archive Using TF*PDF Algorithm," Proc. Int. Conf. on Web Information Systems Engineering (WISE ۲۰۰۲). Singapore, pp. ۷۲-۸۲, ۲۰۰۲.
- [۶] Y. Ko, J. Park and J. Seo, "text categorization using the importance of sentences," Proc. Int. Conf. on Information Processing and Management. South Korea, pp. ۶۵-۷۹, ۲۰۰۴.
- [۷] W. T. Chuang and J. Yang, "Extracting sentences segments for text summarization: a machine learning approaches," Proc. Int.