

۱-۵- اهداف پایان‌نامه

به طور خلاصه اهداف پایان‌نامه را می‌توان به صورت زیر در نظر گرفت:

- مطالعه الگوریتم‌های خلاصه سازی متن
- مطالعه ماشین تورینگ^۱ و کاربردهای آن
- مطالعه اتوماتای یادگیر سلولی^۲ و کاربردهای آن
- مطالعه الگوریتم زنبورعسل مصنوعی^۳ و کاربردهای آن
- مطالعه الگوریتم ژنتیک^۴ و کاربردهای آن
- مطالعه الگوریتم بهینه سازی توده ذرات^۵ برای سیستم امتیازدهی به ویژگی های استخراجی متن
- مطالعه سیستم های استنتاج فازی و کاربردهای آن
- ارائه یک روش جدید برای پیش پردازش متن با ماشین تورینگ
- ارائه یک روش جدید برای تعیین شباهت با استفاده از اتوماتای یادگیر سلولی، الگوریتم زنبورعسل مصنوعی و فرمول شباهت
- ارائه یک روش جدید خلاصه سازی مبتنی بر استخراج ویژگی های مهم متن
- ارائه یک روش جدید خلاصه سازی مبتنی بر استخراج ویژگی های مهم متن، و زدن دهی براساس الگوریتم بهینه سازی توده ذرات
- ارائه یک روش جدید خلاصه سازی مبتنی بر استخراج ویژگی های مهم متن، و زدن دهی براساس الگوریتم بهینه سازی توده ذرات و الگوریتم ژنتیک
- ارائه یک روش جدید خلاصه سازی مبتنی بر استخراج ویژگی های مهم متن، و زدن دهی براساس الگوریتم بهینه سازی توده ذرات و الگوریتم ژنتیک
- ارائه یک روش جدید خلاصه سازی ترکیبی مبتنی بر ۴ روش خلاصه سازی مطرح شده
- مقایسه و ارزیابی روش‌های پیشنهادی با روش‌های دیگر خلاصه سازی

^۱Turing Machine (TM)

^۲Cellular Learning Automata (CLA)

^۳Artificial Bee colony (ABC)

^۴Particle Swarm Optimization (PSO)

۱-۶- ساختار پایان نامه

مطالبی که در پایان نامه حاضر بدان خواهیم پرداخت بدین صورت خواهد بود که در فصل دوم این پایان نامه موضوع خلاصه سازی متن بصورت کلی مورد بررسی قرار می گیرد و انواع رویکردها و تاریخچه ای از آن ها بیان می شود. فصل سوم مفاهیم پایه ای بنام ماشین تورینگ، اتوماتای یادگیر سلولی، الگوریتم بهینه سازی توده ذرات، الگوریتم زنبورعسل مصنوعی، الگوریتم ژنتیک و سیستم های فازی که نقش کلیدی در پایان نامه را دارند، مورد بررسی قرار خواهند گرفت. در فصل چهارم، یک سیستم پیش پردازش پیشنهادی با ماشین تورینگ، یک سیستم تعیین شباهت پیشنهادی با الگوریتم اتوماتای سلولی یادگیر و الگوریتم زنبورعسل مصنوعی و هم چنین ۵ روش خلاصه ساز پیشنهادی به صورت سیستم خلاصه سازی مبتنی بر استخراج ویژگی های متنی، سیستم خلاصه سازی مبتنی بر استخراج ویژگی های متنی و الگوریتم توده ذرات، سیستم خلاصه سازی مبتنی بر استخراج ویژگی متنی، الگوریتم توده ذرات، الگوریتم ژنتیک، سیستم خلاصه سازی مبتنی بر استخراج ویژگی متنی، الگوریتم توده ذرات، سیستم های فازی و سیستم خلاصه سازی ترکیبی مبتنی ۴ روش خلاصه سازی پیشنهادی معرفی می شوند. در فصل پنجم نحوه ارزیابی و پیاده سازی انجام گرفته، همچنین کارایی خلاصه سازهای پیشنهادی در مقایسه با یکدیگر و هم با چنین با خلاصه سازهای پایه H2-H1, Msword, Sys19, Sys30 و خلاصه سازهای دیگر بررسی می شود. در پایان به نتیجه گیری و بیان پیشنهادات می پردازیم.

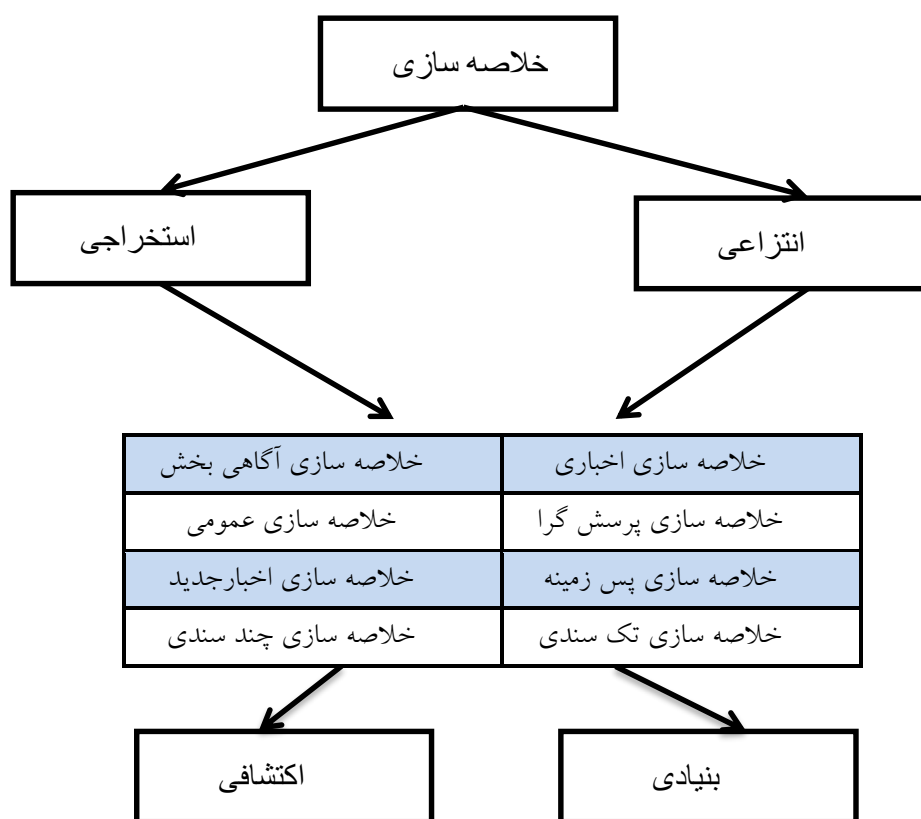
فصل دوم:

متدها و روش‌های خلاصه

سازی متون

۱-۲- مقدمه

خلاصه سازی یک متن را می توان از جنبه های مختلف بررسی کرد و براین اساس خلاصه سازی دارای انواع گوناگونی می باشد. نمایی از انواع خلاصه سازی در شکل ۱-۲ معرفی شده است. حال به معرفی هریک از آنها می پردازیم.



شکل ۱-۲: تقسیم بندی انواع خلاصه سازی

۲-۲- تقسیم بندی انواع خلاصه سازی

۲-۲-۱- خلاص سازی استخراجی در برابر انتزاعی

خلاصه سازی استخراجی، خلاصه‌ای است که عموماً توسط انتخاب بخش‌های مرتبط با موضوع متن بدست می‌آید. خلاصه سازی انتزاعی تفسیری است که محتوای یک مستند را بدون آن که لزوماً از محتوای آن استفاده کرده باشد، شرح می‌دهد. در روش استخراجی، جملات خلاصه معمولاً از حد متوسط طولانی‌تر هستند. اطلاعات مهم معمولاً در طول جملات گسترش یافته‌اند اما جملات استخراج شده این موضوع را در بر نمی‌گیرد مگر اینکه خلاصه به قدر کافی بزرگ باشد که تمام آن جملات را ننگه دارد. در این روش اطلاعات ناسازگار ممکن است به دقت نشان داده نشود. با این حال کاربران خلاصه‌های استخراجی را به چکیده ترجیح می‌دهند. زیرا خلاصه‌های استخراجی، اطلاعات را همان‌طور که نویسنده بیان کرده نشان می‌دهند و پردازش آن‌ها ارزان‌تر است. در روش انتزاعی، خلاصه بوسیله ساختن دوباره عبارات و تشکیل مجموعه‌ای پیوسته که حاوی محتوای مهم متن اصلی است، تهیه می‌گردد. واضح است که در این روش باید از انواع روش‌های گرامری برای ساخت مجدد جمله استفاده نمود که پیچیدگی کار را بالا می‌برد. به همین دلیل این روش امروزه کمتر مورد توجه می‌باشد.

۲-۲-۲- خلاصه سازی اخباری در برابر آگاهی بخش

در خلاصه سازی اخباری، موضوعات اصلی و مهم از متن اصلی استخراج می‌شود که فرد با خواندن این نوع خلاصه، از رئوس اصلی مطالب متن اصلی مطلع می‌شود. ولی برای آگاهی کامل در صورت لزوم باید به اصل متن مراجعه کند. خلاصه سازی آگاهی بخش، خلاصه‌ای کلی بوده است که سعی دارد مطالب مهم را به همراه توضیح کاملشان استخراج نماید به نوعی که نمایانگر کل متن مطلب اصلی می‌باشد و تنها به بیان رئوس مطالب مهم اکتفا نمی‌کند.

۲-۲-۳- خلاصه سازی پس زمینه در برابر اخبار جدید^۲

این دو خلاصه سازی با فرض این که خواننده دانش پیشین در مورد مطلب مورد نظر برای خلاصه را ندارد در مقابل کسی که اطلاعاتش بروز است، می‌باشد.

^۱Background

^۲Just-the-news

۲-۲-۴- خلاصه سازی عمومی در برابر پرسش گرا

خلاصه سازی عمومی در واقع حس کلی محتوای متن را می دهد در حالی که خلاصه های مرتبط با پرس و جو خود را به محتوایی محدود می کنند که با سوال کاربر در ارتباط است. نوع دوم به طور وسیع در ارتباط با مستندات که بزرگ هستند و از لحاظ موضوع متفاوتند بسیار کار آمد تر است.

۲-۲-۵- خلاصه سازی تک سندی^۱ در برابر چند سندی^۲

در حالت تک سندی خلاصه تهیه شده از روی یک متن بوده و در خلاصه چند سندی خلاصه تولیدی بوسیله ترکیب چند متن و استخراج خلاصه می باشد.

۲-۲-۶- خلاصه سازی براساس دیدگاه های اکتشافی^۳ و بنیادی^۴

به طور کلی متدها و روش های خلاصه سازی را می توان به دودیدگاه کلی تقسیم بندی نمود: دیدگاه های اکتشافی و دیدگاه های بنیادی. تکنیک هایی که براساس نظریات شناخته شده ریاضی عمل می کنند، در دسته دیدگاه های بنیادی تقسیم بندی می شوند. تکنیک هایی که از تجربیات انسان استفاده می کنند در دسته دیدگاه های اکتشافی قرار می گیرند، از این دسته به تحلیل های تجربی زبان طبیعی و تحلیل نحوه نوشتار انسان اشاره کرد. در اکثر مواقع دودیدگاه برای دستیابی به خلاصه سازی بهتر به کمک یکدیگر می آیند. در ادامه برخی از این روش هایی که بر پایه این دو دیدگاه می باشند را توضیح می دهیم.

۲-۲-۳- انواع روش های خلاصه سازی

۲-۳-۱- روش بر پایه موقعیت^۵

این روش برای اولین بار توسط Edmundson در سال ۱۹۶۹ اعمال گردید و در دسته تکنیک های اکتشافی گنجانده شده است. ایده این روش بدین صورت است که جملات مهم متن در ابتدا یا انتها قرار دارند. تجربه نشان داده که در ۸۵ درصد موارد جملات مهم در مکان های ابتدایی و در ۷ درصد موارد در مکان های پایانی قرار دارند. همچنین تنها ۱۳ درصد از پاراگراف های نویسنده های معاصر با جملات مهم شروع می شوند (Baxendale, 1958). نقطه ضعف این روش این است که محل جملات مهم در پاراگراف

¹Single Document

²Multi Document

³Heuristic approach

⁴Foundation-based approach

⁵Position-Based Method

می تواند متناسب با موضوع متون متفاوت باشد. به عبارتی می توان موقعیت جمله مهم را در متن با استفاده از یک فرآیند یادگیری بدست آورد. نقطه ضعف دیگر این روش این است که برای کلیه پاراگراف های متن ارزش یکسانی در نظر می گیرد. به عبارتی ممکن است پاراگراف های ابتدایی یا انتهایی از اهمیت بیشتری برخوردار باشند.

۲-۲-۳-۲- روش سیاست بهینه مکانی^۱

این روش اولین بار توسط Lin و Hovey در سال ۱۹۹۷ پیشنهاد گردید (Lin and Hovy, ۱۹۹۷). این روش در دسته تکنیک های اکتشافی گنجانده شده است. ایده این روش بدین صورت است که جملات مهم در متون مختلف با نوع موضوع در مکان مختلفی واقع می شوند و با استفاده از یک الگوریتم یادگیری این مکان به صورت خودکار تعیین می گردد. در این روش برای هر متن همپوشانی بین هر یک از جملات متن با کلمات کلیدی و جملات چکیده را تعیین می نماییم. براساس این همپوشانی و شباهت یک امتیاز به جمله مذکور تعلق می گیرد. در قدم بعدی با دو مسئله روبرو هستیم:

۱. برای هر پاراگراف مجموع میانگین امتیازات جملات آنرا حساب می کنیم. بدین ترتیب می توان پاراگرافی که دارای بیشترین اهمیت است را استخراج نمود.

۲. برای هر موقعیت جمله در پاراگراف، برای مثال جملات اول پاراگراف های امتیاز را حساب می کنیم. بنابراین مشخص می شود که به طور متوسط در پاراگراف های مختلف چندیم جمله از بیشترین امتیاز برخوردار است.

بنابر دو مسئله گفته شده در بالا، جمله مهم قابل استخراج بوده یا به عبارتی جمله ای که در شرایط بالا صدق می کند از امتیاز بیشتری برخوردار است.

مجموعه داده مورد استفاده در این روش مجموعه داده ZIFF بوده که برای آموزش ۱۳۰۰۰ مقاله روزنامه در مورد تازه های کامپیوتر و سخت افزار و غیره می باشد. نوع و دسته بندی هر یک از این مقالات نیز در این مجموعه داده مشخص گردیده است. هر مقاله بطور متوسط شامل ۷۱ جمله بوده و همراه با هر مقاله ۳ تا کلمه کلیدی و یک خلاصه که شامل ۶ جمله تولید شده بوسیله انسان می باشد. نتایج بدست آمده برای ۲۹۰۰ مقاله از این مجموعه مطابق با جدول ۱-۲ بوده است.

جدول ۱-۲: نتایج بدست آمده از روش سیاست بهینه مکانی

دقت	فراخوانی
-----	----------

^۱ Optimum Position Policy

۳۵٪	۳۸٪
-----	-----

این روش درمقایسه با روش برپایه موقعیت بهتر عمل کرده است زیرا این روش با توجه به اینکه موقعیت مناسب جملات و پاراگراف را به صورت داینامیکی و متناسب با نوع متن بدست می آوردند از کارایی بهتری برخوردارست. در روش برپایه موقعیت، فرض بر این بود که جملات مهم در پاراگراف اول می باشد ولی این فرض در حالت کلی درست و کارا نمی باشد.

۲-۲-۳- روش برپایه عنوان

این روش نیز برای اولین بار توسط Edmundson در سال ۱۹۶۹ اعمال گردید و در دسته تکنیک های اکتشافی گنجانده شده است. ایده این روش بدین صورت است که کلمات واقع شده در عنوان های کلی و عنوان های بخش های جزئی تر متن ارتباط مستقیمی با خلاصه تولید شده دارند. تجربه نشان داده که ۹۹٪ این جملات به صورت آماری از اهمیت برخوردار هستند و در خلاصه کردن مفید واقع می شوند.

۲-۲-۳-۴- روش برپایه عبارات خاص^۱

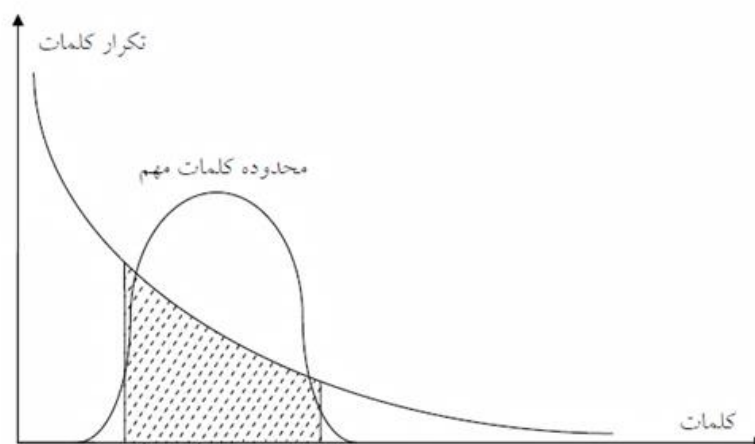
این روش در دسته تکنیک های اکتشافی گنجانده می شود، نیز برای اولین بار توسط Edmundson در سال ۱۹۶۹ اعمال گردید. جملات مهم متن اغلب شامل عبارات امتیازی مانند *In*، *Significantly*، *In conclusion*، *We show*، *this paper* و غیره می باشند. در مقابل این عبارت، عبارات با امتیاز منفی مانند *Hardly* یا *impossible* می باشد. این عبارات نیز خود می توانند به صورت خودکار از متن استخراج شوند. نحوه کارچنین است که با توجه شدن به یک عبارت با امتیازی یک امتیاز مثبت و با توجه شدن به یک عبارت با امتیاز منفی یک پنالتهی به جمله مورد نظر اضافه می کنیم و در نهایت جمله با امتیاز بیشتر را استخراج می کنیم.

نقطه ضعف این روش در این است که عبارات خاص باید از پیش تعیین شوند که این خود متوجه هزینه و در ضمن اینکه عبارات خاص در متون مختلف متفاوت هستند که برای حل این مشکل Tufel در سال ۹۸ روشی برای استخراج خودکار این عبارات ارائه نموده است.

^۱Cue-Based Method

۲-۲-۳-۵- روش برپایه بسامد لغوی^۱

این روش که در دسته تکنیک های اکتشافی گنجانده می شود و بوسیله Luhn در سال ۱۹۵۸ معرفی گردید و بدین صورت بود که جملات مهم حاوی کلماتی هستند که دارای بسامد تکرار متوسطی می باشند. یعنی نه تکرار آنها بسیار زیاد باشد مانند حروف تعریف یا ضمائر و... یا بسیار نادر باشند. بدین منظور جملات را به ازای هر کلمه با ویژگی گفته شده در آن برامتیازش می افزاییم. نمودار بسامد تکرار کلمات در شکل ۲-۲ آورده شده است.



شکل ۲-۱. نمودار بسامد تکرار کلمات

در این روش برای شمارش کلمات مشابه از ریشه^۲ یابی نیز استفاده می شود. منظور از ریشه یابی این است که برای مثال دو کلمه system و systems مشابه هستند. در این روش همچنین از یکسری کلمات مانند is ، a ، the و... صرف نظر می شود و Luhn از آنها با نام لغات توقف^۳ یاد کرده است. در قدم بعد باید tf هر کلمه در متن محاسبه گردد. منظور از tf محاسبه توزیع احتمالی وقوع یک کلمه در متن هدف می باشد. همچنین باید idf هر کلمه نیز محاسبه شود. مقدار idf یک عبارت از معادله ۱-۲ محاسبه می گردد:

$$idf(term) = \log\left(\frac{NUMDOC}{NUMDOC(term)}\right) \quad (1-2)$$

که در آن NOMDOC تعداد متن های موجود در پیکره بوده و NUMDOX(term) تعداد متن هایی است که عبارت term در آنها رخ داده است. در نهایت الگوریتم بدین صورت است که عبارتی

^۱Word-Frequency-Based Method

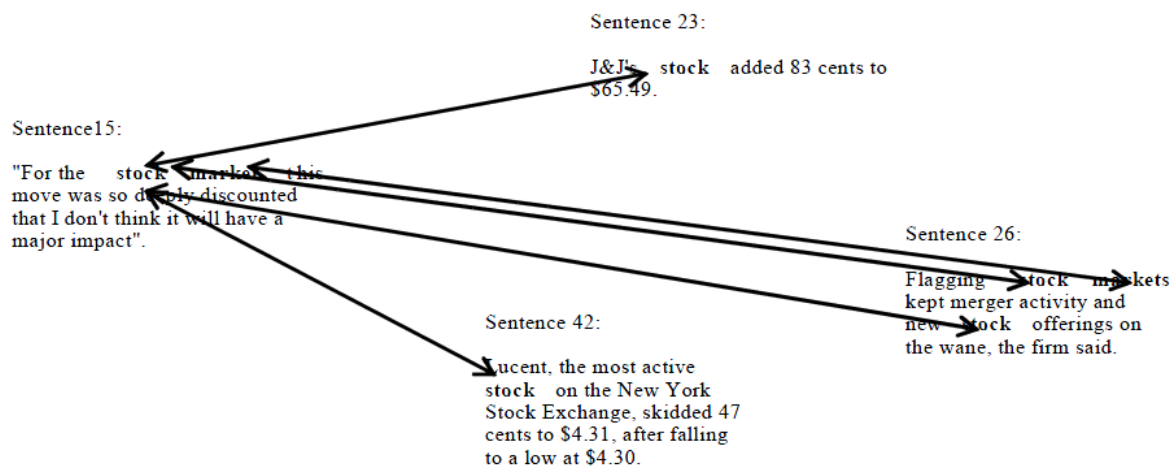
^۲Stem

^۳Stop words

انتخاب می گردد که مقدار $tf \times idf$ آن بیشتر از یک حد آستانه باشد. بنابراین می توان وزن اختصاصی به یک جمله را از مجموع $tf \times idf$ های لغات آن و نرمال کردن آنها بدست آورد.

۲-۳-۲-۶- روش های برپایه پیوستگی^۱

در ابتدا بوسيله Ahmad و Bembrahim در سال ۱۹۹۵ معرفی گردید و این روش در دسته تکنیک های بنیادی گنجانده شده است. ایده این روش ها بدین صورت است که جملات و عبارات مهم دارای بیشترین اتصال درگراف شباهت خود هستند. این گراف خود می تواند براساس تکرار یک کلمه یا شباهت معنایی یا هم مرجعی و غیره باشد (Ponte and Croft, 1998). در این روش چندین دیدگاه وجود دارد که هر یک را بصورت مجزا مورد بررسی قرار خواهیم داد. شکل زیر یک نمونه ارتباط بین چند جمله را آورده شده است. همانطور در شکل ۳-۲ می بینید، می خواهیم ارتباط ۴ جمله فرضی را براساس هم وقوعی کلمات آن بدست آوریم. برای مثال دو جمله ۱۵ و ۲۳ در کلمه stock مشترک هستند، به عبارتی کلمه stock در این جمله هم وقوع شده است. برای نشان دادن این ارتباط یک یال بین این دو جمله رسم می کنیم. برای مثالی دیگر دو جمله ۲۶ و ۱۵ در دو کلمه stock و market مشترک می باشند. بنابراین گراف مورد نظر حاصل می شود.



شکل ۳-۲: یک نمونه ارتباطی لغوی در چند جمله

از جمله روش هایی که از این متد استفاده کرده اند به صورت زیر می باشد:

- روش هم وقوعی کلمات^۱

• شباهت لغوی (زنجیره لغوی^۲، WorldNet)

• ترکیبی از روش های بالا

۲-۲-۳-۷- روش هم وقوعی کلمات

این روش بنیادی برای اولین بار بوسیله Mitra در سال ۱۹۹۷ ارائه گردید. در این روش از تکنیک های بازیابی اطلاعات در سطح متن استفاده می شود. به عبارتی متون را می توان مجموعه ای از پاراگراف ها در نظر گرفت. با استفاده از متدهای پیشین در بازیابی اطلاعات و با استفاده از معیارهای شباهت بین کلمات، شباهت بین پاراگراف p_i را با دیگر پاراگراف های متن تعیین می کنیم (Salton, Singhal, Mitra and Buckley, 1997). این روش بدین صورت است که ابتدا باید یک واحد برای متن خود جهت تشابه انتخاب کنیم که Mitra در مقاله خود این واحد را پاراگراف انتخاب کرده است. در قدم بعدی هر پاراگراف را به صورت یک بردار در نظر می گیریم که عناصر آن وزن یا اهمیت ترم (d_{ik}) مورد نظر در پاراگراف است که این ترم می تواند کلمه باشد. در قدم بعدی باید بتوانیم شباهت بین دو پاراگراف را مشخص کنیم. با در اختیار داشتن بردار ۲-۲

$$D_i = (d_{i1}, \dots, d_{im}) \quad (۱-۲)$$

در قدم بعدی باید شباهت بین پاراگراف را مشخص کنیم. با در اختیار داشتن بردار D_i و نرمال کردن آن شباهت دو پاراگراف از فرمول ۲-۳ بدست می آید:

$$\text{Sim}(D_i, D_j) = \sum d_{ik} d_{jk} \quad (۳-۲)$$

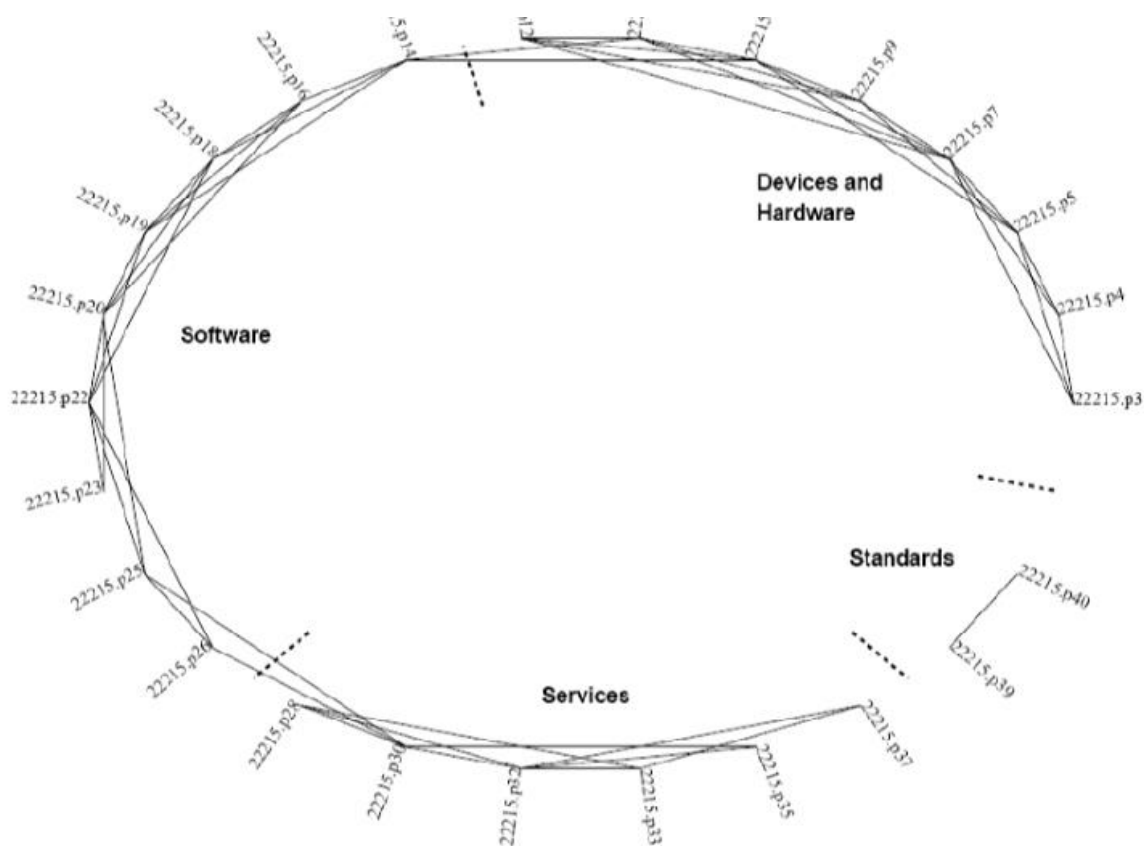
سپس کلیه پاراگراف را به هم متصل نموده و وزن اتصالات را برابر مقدار شباهت آنها قرار می دهیم. بدین ترتیب برای یک متن به گرافی بدست می آید که به آن نقشه ارتباطی متن گفته می شود. در قدم بعد باید یال هایی که وزن کمتر از بیک مقدار آستانه دارند را حذف کنیم که در نهایت به گرافی مانند شکل زیر می رسیم.

با دقت در گراف بدست آمده در شکل ۴-۲ متوجه می شویم که برای مثال دو پاراگراف p39 و p40 میزان ارتباطی بیشتر از حد آستانه را نداشته اند. بنابراین تشکیل یک حوزه یا خوشه معنایی را می دهند که در این مثال در بردارنده معنای Standards هستند. بنابراین گراف ما با حذف ارتباطات کمتر از مقدار آستانه، به چند بخش تقسیم می شوند که اتصالات درون هر بخش بیشینه و اتصال بین بخش ها کمینه

^۱WORD co-occurrence

^۲Lexical chains-based method

است. این زیر گراف ها در واقع پاراگراف هایی را نشان می دهند که از لحاظ معنایی بیشترین ارتباط را دارند و همگی تقریباً یک مفهوم را دربردارند. برای مثال در گراف مثال قبل این بخش ها نشان داده شده اند.



شکل ۴-۲: یک نمونه ارتباط بین پاراگراف ها در یک متن

در قدم بعدی باید از بین پاراگراف های یک بخش، مهمترین آنها انتخاب کنیم. برای این انتخاب روش های متنوعی ارائه گردیده که برای مثال می توان پاراگرافی را انتخاب نمود که بیشترین اتصال را داشته باشد. برای ارزیابی نیز بدین طریق عمل شده که ابتدا پاراگراف های مهم یک متن توسط دو فرد استخراج گردیده و سپس میزان همپوشانی پاراگراف های استخراجی این دو فرد مشخص می گردد. سپس همپوشانی پاراگراف های استخراج شده بوسیله سیستم را در بهترین حالت و بدترین حالت و اشتراک پاراگراف های