



بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه اصفهان  
دانشکده فنی و مهندسی  
گروه مهندسی کامپیوتر

## پایان نامه کارشناسی ارشد رشته‌ی مهندسی کامپیوتر گرایش نرم افزار

خلاصه‌سازی خودکار متون فارسی به روش معنایی

استادان راهنما:

دکتر محمد علی نعمت بخش

دکتر احمد رضا نقش نیلچی

پژوهشگر:

فراز محمدیان جدول قدم

شهریورماه 1391

کلیه حقوق مادی مترتب بر نتایج مطالعات، ابتکارات  
و نوآوری‌های ناشی از تحقیق موضوع این پایان‌نامه  
متعلق به دانشگاه اصفهان است.



دانشگاه اصفهان  
دانشکده فنی و مهندسی  
گروه کامپیوتر

پایان نامه‌ی کارشناسی ارشد مهندسی کامپیوتر گرایش نرم افزار  
آقای فراز محمدیان جدول قدم تحت عنوان

**ارائه یک روش خلاصه‌ساز خودکار متون فارسی به روش معنایی**

در تاریخ 91/12/13 توسط هیأت داوران زیر بررسی و با درجه بسیار خوب به تصویب نهایی رسید.

- |      |   |                                     |
|------|---|-------------------------------------|
| امضا | دکتر محمدعلی نعمت بخش با مرتبه‌ی علمی دانشیار | 1- استاد/استادان راهنمای پایان نامه |
| امضا | دکتر احمد نقش نیلچی با مرتبه‌ی علمی دانشیار   | 2- استاد/استادان راهنمای پایان نامه |
| امضا | دکتر بهمن زمانی با مرتبه‌ی علمی استادیار      | 3- استاد/استادان داور داخل گروه     |
| امضا | دکتر بتول علی نژاد با مرتبه‌ی علمی دانشیار    | 4- استاد/استادان داور خارج از گروه  |

امضای مدیرگروه

## با تشکر و قدردانی از

اساتید گرانقدر آقایان دکتر محمدعلی نعمت بخش و دکتر احمد رضا نقش نیلچی  
که اساتید راهنمای این پروژه بودند و مرا در تمام مراحل انجام این کار یاری  
رساندند.

تقدیم به

پدر و مادر و همسر

## چکیده:

خلاصه‌سازی خودکار متن یکی از جنبه‌های مهم در بازیابی اطلاعات می‌باشد. در این پایان‌نامه یک روش جدید برای خلاصه‌سازی خودکار متون فارسی پیشنهاد شده است که رویکردی مبتنی بر معنا دارد. روش پیشنهادی، دارای سه فاز پیش پردازش، پردازش و تولید خلاصه نهایی می‌باشد. در مرحله پیش پردازش، ریشه کلمات با استفاده از یک روش مبتنی بر فرهنگ واژگان بدست می‌آید. مجموعه مترادف جملات موجود در متن از طریق همین فرهنگ واژگان بدست می‌آید. این کار منجر به یک زنجیره لغوی قوی می‌شود. در مرحله پردازش، با استفاده از زنجیره لغوی و روش تحلیل معنای پنهان، ارتباط بین جملات مهم بدست می‌آید تا جملات مرتبط با هم گزینش و در خروجی قرار گیرند. در مرحله تولید خلاصه، برای رفع افزونگی و تکرار، جملات مشابه در خروجی با جملات انتخاب شده مقایسه شده و در صورت داشتن شباهت نسبی (با انتخاب ضریب میزان تشابه 80٪) از قرار دادن آن در خلاصه نهایی جلوگیری می‌شود.

دو شیوه برای ارزیابی خلاصه‌ساز پیشنهادی ارائه شده است. یکی روش ارزیابی مستقیم و دیگری روش ارزیابی مبتنی بر LSA (تحلیل معنایی پنهان) است. نتایج ارزیابی نشان می‌دهد به دلیل استفاده از فرهنگ واژگان غنی، ریشه‌یابی با دقت بالایی انجام می‌گیرد که در مقایسه با روش‌های مشابه کارایی بالایی دارد. به دلیل استفاده از همین فرهنگ واژگان غنی، مجموعه مترادف و زنجیره لغوی ارتباط معنای کلمات و جملات در روش تحلیل معنایی پنهان بسط داده شده به درستی بدست آورده می‌شود که نتایج حاکی از این امر دارد. در نهایت نتایج ارزیابی نشان می‌دهد که روش ارائه شده خلاصه قابل قبول و منسجمی ارائه می‌دهد.

**واژگان کلیدی:** خلاصه‌سازی، ریشه‌یابی، زنجیره لغوی، تحلیل معنایی پنهان، ارتباط معنایی.



## فهرست مطالب

صفحه

عنوان

### فصل اول: طرح تحقیق

1-1-1- مقدمه ..... 1

### فصل 2: خلاصه‌سازی و خلاصه‌سازی خودکار

1-2-1- مقدمه ..... 5

2-2- خلاصه‌سازی ..... 5

3-2- خلاصه‌سازی خودکار ..... 6

4-2- فرآیند خلاصه‌سازی ..... 7

1-4-2- پردازش لغوی ..... 7

2-4-2- پردازش ساخت واژگی ..... 8

5-2- طبقه‌بندی سیستم‌های خلاصه‌سازی ..... 9

6-2- مراحل خلاصه‌سازی ..... 11

1-6-2- پیش پردازش متن ..... 11

2-6-2- پردازش متن خلاصه استخراجی ..... 14

1-2-6-2- روش سطحی ..... 14

2-2-6-2- روش موجودیتی - معنایی ..... 14

3-2-6-2- سطح کلامی ..... 15

5-2-6-2- روش ترکیبی ..... 17

7-2- پردازش متن خلاصه چکیده ..... 17

8-2- تولید خلاصه ..... 17

### فصل 3: کارهای انجام شده

1-3-1- مقدمه ..... 19

1-1-3- خلاصه ساز تجربی ..... 19

2-1-3- گرامر موردی ..... 19

2-3- خلاصه استخراجی ..... 20

1-2-3- روشهای بر پایه پردازش سطحی ..... 20

2-2-3- روشهای موجودیتی - بر پایه معنی ..... 24

3-2-3- روشهای بر پایه ساختار کلامی ..... 34

37	4-2-3- روشهای ترکیبی
43	3-3- خلاصه‌سازی چند سنده
44	4-3- خلاصه‌سازهای کارشده در زمینه فارسی

#### فصل 4: روش پیشنهادی

45	1-4- مقدمه
45	2-4- روش پیشنهادی
46	3-4- پیش‌پردازش
46	1-3-4- یکسان‌سازی دستور خط فارسی
50	2-3-4- شناسایی ریشه فعل‌ها
51	3-3-4- تعیین مرز کلمات املائی
52	4-3-4- حذف واژه‌های عمومی
55	5-3-4- ریشه‌یابی کلمات
63	6-3-4- پیدا کردن مترادف کلمات در فرهنگ واژگان (ایجاد زنجیره لغوی)
64	4-4- مرحله پردازش متن
64	1-4-4- زنجیره لغوی
65	2-4-4- مرحله تحلیل معنایی
70	5-4- تولید خلاصه
71	6-4- نتیجه‌گیری

#### فصل پنجم: ارزیابی و نتایج تجربی

72	1-5- مقدمه
72	2-5- روش ارزیابی
74	1-2-5- ارزیابی مستقیم
76	2-2-5- ارزیابی سامانه‌های مشابه و پیشنهادی
78	3-2-5- ارزیابی مبتنی بر LSA
80	3-5- نتایج علمی روش ارزیابی مبتنی بر LSA
82	4-5- نتیجه‌گیری

عنوان

صفحه

فصل ششم: جمع بندی و کارهای آینده

83	1-6- مقدمه
83	2-6- نتیجه گیری
84	3-6- کارهای آینده
85	منابع

## فهرست شکل ها

عنوان	صفحه
شکل 1-1- فرآیند خلاصه‌سازی .....	2
شکل 2-2- خلاصه‌سازی خودکار سند .....	6
شکل 1-3- درخت ساختار زبانی .....	35
شکل 2-3- مثالی از درخت ساختار کلامی .....	35
شکل 1-4- معماری روش پیشنهادی .....	46
شکل 2-4- حذف پیشوندها و پسوندها .....	62
شکل 3-4- الگوریتم ریشه‌یابی لغات با استفاده فرهنگ واژگان دهخدا .....	62
شکل 4-4- شبه کد الگوریتم ریشه‌یابی پیشنهادی .....	63
شکل 5-4- الگوریتم بدست آوردن زنجیره لغوی .....	65
شکل 4-6- الگوی ساده کلمه - مفهوم .....	66
شکل 4-7- الگوی مختلف کلمات و معانی به هم مرتبط .....	66
شکل 4-8- مقدار تجزیه منفرد (SVD) .....	67
شکل 4-9- ماتریس کلمه - جمله (زنجیره لغوی) .....	69
شکل 4-10- الگوریتم پردازش متن و تولید خلاصه .....	70
شکل 1-5- نتایج مربوط به مستند 1 .....	77
شکل 2-5- نتایج مربوط به مستند 2 .....	77
شکل 3-5- میانگین اولویت‌بندی‌ها .....	78
شکل 4-5- اولین بردار منفرد (عنوان اصلی) .....	79
شکل 5-5- ایجاد بردار واژگان حاصل از یک متن اصلی و خلاصه .....	80

## فهرست جداول

صفحه	عنوان
32	جدول 3-1- نتایج ارزیابی
51	جدول 4-1- هشت گروه فعلهای فارسی
52	جدول 4-2- تعیین مرز املائی کلمات
52	جدول 4-3- الگویی صحیح تعیین مرز املائی کلمات
53	جدول 4-4- تمام اشکال مختلف فعل "کردن"
54	جدول 4-5- حروف اضافه فعلی
54	جدول 4-6- حروف اضافه غیر فعلی
75	جدول 5-1- مقایسه نتایج خلاصه مرجع و خلاصه سامانه پیشنهادی مستند 1
76	جدول 5-2- مقایسه نتایج خلاصه مرجع و خلاصه سامانه پیشنهادی مستند 2
81	جدول 5-3- ارزیابی شباهت کسینوسی - ارزیابی مبتنی بر محتوی
81	جدول 5-4- ارزیابی شباهت عنوان اصلی
81	جدول 5-5- ارزیابی شباهت واژگان مهم

## فصل اول

### طرح تحقیق

#### 1-1- مقدمه

امروزه بیش از 80 درصد از دانش ما به صورت متن، مستندات<sup>1</sup> و دیگر صورت‌های رسانه‌ای نظیر ویدئو و صدا نگهداری می‌شود. اگر از منظر علوم کامپیوتری به این مستندات نگاه کنیم همه‌ی آنها به طبیعتی غیرساخت‌یافته وابسته‌اند. با رشد روز افزون منابع اطلاعاتی، پید کردن منابع مربوط به موضوع مورد نظر و گزینش<sup>2</sup> مطالب مفید از میان منابع مربوطه به یک معضل تبدیل شده است. این مشکل توسط سیستم‌های خلاصه‌سازی قابل حل است. از این رو پیدا کردن اطلاعات متنی مناسب و ارائه یک چکیده<sup>3</sup> (خلاصه) به یک چالش بزرگ تبدیل شده است [1].

#### 1-2- خلاصه

خلاصه، در لغت به معنی خالص، برگزیده، منتخب و کوتاه‌شدن مطلب است. خلاصه، یک ارائه دقیق از محتوای یک متن است. برای رسیدن به این منظور در زیر تعریفی از خلاصه ارائه شده است: «کاهش و تغییرات متن اصلی به خلاصه متن، از طریق انتخاب یا عمومیت دادن به آنچه که در متن است.» استفاده از خلاصه، زمان خواندن متن را متناسب با اطلاعات مورد نیاز تا حد زیادی کاهش می‌دهد. یک فرد برای دریافت دانش از اطلاعات یک متن، بایستی ابتدا آنرا درک کند، سپس آن را پردازش کرده تا بفهمد چه معانی و مفاهیمی در متن موجود است و از میان این مفاهیم کدام جدید و کدام قدیمی است. بعد از درک مفاهیم، خلاصه‌ای متناسب با

---

1. Document's  
2. Election  
3. Abstract

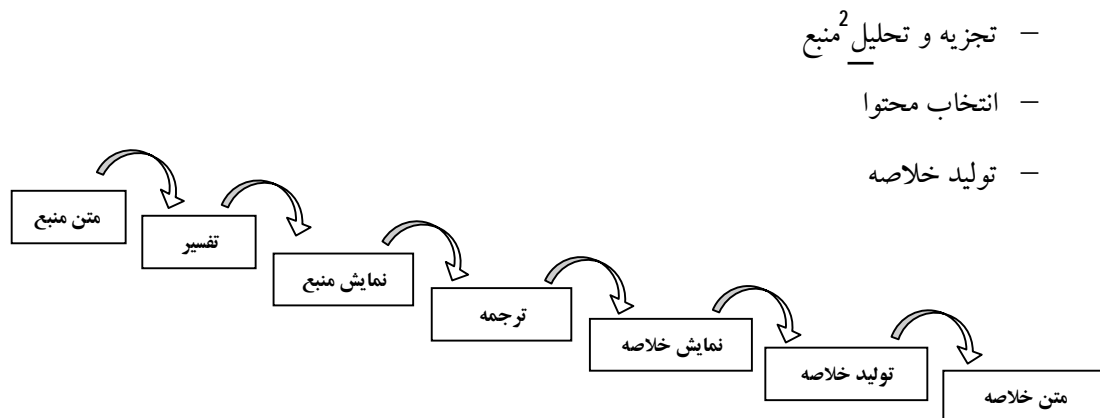
نیاز خود تولید نماید. این کار سخت و طاقت فرسا می‌باشد و از طرفی چون قضاوت انسانی بی‌ثبات است و ماهیت افکار آن‌ها در محیط‌های مختلف قابل تغییر است، خلاصه‌های مختلفی را می‌توان از افراد داشته باشیم. بنابراین اگر بتوان خلاصه را توسط ماشین و بصورت خودکار تولید نمود این مشکل قابل حل است. این رهیافت مستلزم استفاده از خلاصه‌سازی خودکار متن است [3، 5].

خلاصه‌سازی خودکار متن<sup>1</sup> عبارت است از: کوتاه کردن متن از طریق انتخاب جملات مهم، به طوری که متن نهایی مهم‌ترین نکات متن اصلی را نیز در برداشته باشد. خلاصه‌سازی خودکار متن، تکنیکی است که از طریق یک برنامه کامپیوتری متنی را خلاصه می‌کند و در اختیار درخواست کننده قرار می‌دهد. این برنامه یک متن از ورودی دریافت کرده و یک خلاصه از متن اصلی به خروجی می‌فرستد. بنابراین در خلاصه‌سازی خودکار، اطلاعات از منبع اصلی استخراج، آنگاه محتوای آن بررسی شده و محتوای مهم برای کاربر به شکل متراکم و به هم پیوسته ارائه می‌شود [17، 6، 2].

از جمله مزیت‌های خلاصه‌سازی خودکار متن نسبت به خلاصه‌سازی دستی می‌توان به موارد زیر اشاره کرد:

- ارتباط بین موجودیت‌های متن، در خلاصه و موقعیت آن‌ها در متن اصلی به آسانی قابل ایجاد است.
- اندازه خلاصه و محتوای آن قابل کنترل و معین است.

فرآیند تبدیل منبع به خلاصه در سه مرحله ذیل انجام می‌پذیرد که در شکل 1-1 نشان داده شده اند:



شکل 1-1- فرآیند خلاصه‌سازی

1. Automatic summarization
2. Analysis
3. Preprocessing

فرایند تجزیه و تحلیل منبع، که به آن پیش‌پردازندش<sup>1</sup> هم گفته می‌شود خصیصه‌هائی که به یک خلاصه منتج می‌شود را تعیین کرده و یک نمایش ساخت یافته از متن اصلی به دست می‌آورد. فرایند انتخاب محتوا یا پردازش الگوریتمی، اطلاعات متن را براساس خصیصه‌هائی که در مرحله قبل بدست آمده است به ساختار خلاصه تبدیل می‌کند، نهایتاً در فرایند تولید خلاصه، خلاصه نهایی با استفاده از ساختار خلاصه بدست می‌آید [19، ۳۸].

روش‌های خلاصه‌سازی می‌تواند براساس سطح زبان به دو گروه عمده تقسیم شود:

1. دیدگاه سطحی که محدود به نمایش قواعد نحوی است و سعی در استخراج قسمت‌های برجسته متن به روش مناسبی دارد.

2. دیدگاه عمیق‌تر که شامل سطح معنایی از متن اصلی و فرایندهای زبانی است.

در دیدگاه اول هدف از پیش‌پردازش متن، کاهش بعد فضایی آن است و به صورت طبیعی شامل موارد زیر می‌باشد:

(الف) حذف حروف اضافه<sup>2</sup> (لغات ارتباط دهنده رایج بین کلمات بدون هیچ معنایی خاص مثل a و The)  
(ب) نادیده گرفتن حروف بزرگ و کوچک و تبدیل تمامی کاراکترهای متن به فرم خاص از حروف (حروف بزرگ یا کوچک)

(ج) تبدیل کلمات به ریشه<sup>3</sup> آنها، کلماتی با قواعد نحوی مشابه مثل کلمات جمع، حالت‌های مختلف فعل و غیره به صورت یکسان در نظر گرفته می‌شوند. هدف از این مرحله به دست آوردن ریشه کلمه است که در پردازش معنایی استفاده می‌شود [11، ۵].

### 3-1- انواع خلاصه‌سازی

خلاصه را می‌توان به دو دسته استخراجی<sup>4</sup> و غیراستخراجی (چکیده<sup>5</sup>) تقسیم کرد. در خلاصه‌ی استخراجی جمله‌ها از متن اصلی گزینش می‌شوند و در خلاصه کپی می‌گردند، به عنوان مثال عبارات، جملات یا پاراگراف‌های کلیدی در متن اصلی عیناً در خلاصه کپی می‌شوند. در این روش خلاصه نهایی انسجام و یکپارچگی را بین پاراگراف‌ها و جملات تضمین نمی‌کند ولی نمایش تقریبی از محتوای متن را به نمایش می‌گذارد. خلاصه غیر استخراجی مطالب مفید متن اصلی را، انتخاب می‌کند و آنها را تفسیر می‌نماید و در

---

2. Stop word  
3. Stem  
4. Extractive  
5. Abstracting



خلاصه بازنویسی می‌کند. به طور کلی خلاصه غیراستخراجی قوی‌تر از خلاصه استخراجی می‌تواند متن را خلاصه کند، برنامه‌هایی که این نوع خلاصه را تولید می‌کنند، نیاز به بکارگیری تکنولوژی تولید زبان طبیعی<sup>1</sup> دارند که خود بحث پیچیده‌ای است [17، 2، 1].

خلاصه نیز می‌تواند انواع دیگری داشته باشد که می‌توان به موارد زیر اشاره کرد [2، 1]:

- بر حسب تعداد اسناد ورودی: تک سندی یا چند سندی<sup>2</sup> می‌تواند باشد.
  - بر حسب تعامل با کاربر: مبتنی بر پرس و جو یا غیر مبتنی بر پرس و جو می‌تواند باشد.
  - بر حسب عمومیت: عمومی یا محدود به حوزه می‌تواند باشد.
- خلاصه‌سازی خودکار متن دارای ضعف‌ها و مشکلاتی است. در زیر این چالش‌ها را دسته‌بندی کرده و مورد بررسی قرار خواهیم داد:

- خلاصه‌سازی درگیر مکانیزم تهیه خروجی موثر و کارآمد است. زیرا ایجاد یک خلاصه یکپارچه و منسجم که عبارات آن دارای یک رابطه منطقی با یکدیگر باشند بسیار حائز اهمیت می‌باشد.

- ارزیابی و قضاوت کردن در مورد خلاصه تولیدشده را می‌توان از دیگر چالش‌های خلاصه‌سازی خودکار نام برد. زیرا رهیافت‌های بکار رفته در ارزیابی این سیستم‌ها مانند قضاوت انسانی بی‌ثبات است. در این رهیافت نتایج وابستگی زیادی به افراد دارد و دیدگاه افراد با هم متفاوت است.

- خلاصه‌سازی درگیر انسجام و تراکم<sup>3</sup> بین عبارات موجود در خلاصه است، زیرا توانایی ارزیابی خلاصه‌ها با نرخ تراکم متفاوت، مقیاس و پیچیدگی ارزیابی را افزایش می‌دهد.

- از سوی دیگر در خلاصه‌سازی غیر استخراجی، خلاصه‌ی تولید شده ممکن است بیرون از محتوای متن باشد، به گونه‌ای که تفسیری که خلاصه‌ساز از متن اصلی داشته است اشتباه بوده و هیچ‌گونه ارتباطی از لحاظ محتوا با متن اصلی نداشته باشد. در این حالت مسئله‌ی وابستگی یا انسجام<sup>4</sup> خلاصه تولید شده، نسبت به منبع اصلی بوجود می‌آید [2].

نحوه نمایش خلاصه توسط خلاصه‌ساز از دیگر چالش‌های خلاصه‌سازی خودکار متون است، زیرا باید نحوه نمایش به گونه‌ای پاسخگوی نیازهای افراد در زمینه‌ها و حوزه‌های مختلف باشد.

---

1. Natural language generation  
 2. Multi document  
 3. Compression  
 4. Coherence

## فصل دوم

### خلاصه‌سازی و خلاصه‌سازی خودکار

#### 2-1- مقدمه

در این فصل، ابتدا مفاهیم و تعاریف اولیه مربوط به مبحث بازیابی اطلاعات و خلاصه‌سازی متون مطرح می‌شود. سپس انواع منابع و روش‌های پیشین و مشکلات این روش‌ها را مورد بررسی قرار می‌دهیم.

#### 2-2- خلاصه‌سازی

خلاصه در لغت یعنی خاص، برگزیده، منتخب، کوتاه‌کردن مطلب و ... مصدر آن معنی خلاصه کردن، مختصر کردن کلام و روشن ساختن می‌باشد. خلاصه در مقابل واژه چکیده نیز آمده است و چکیده یعنی نوشته‌ای از مهمترین مطالب یک کتاب و یا یک مقاله و متن گزارش که بصورت خلاصه تهیه شده است (شکل 1-2). منظور از خلاصه نمودن در اینجا بیان اصل مطلب در کوتاه‌ترین شکل می‌باشد، بطوریکه به هیچ وجه به اصل مطلب خدشه‌ای وارد نیاید.

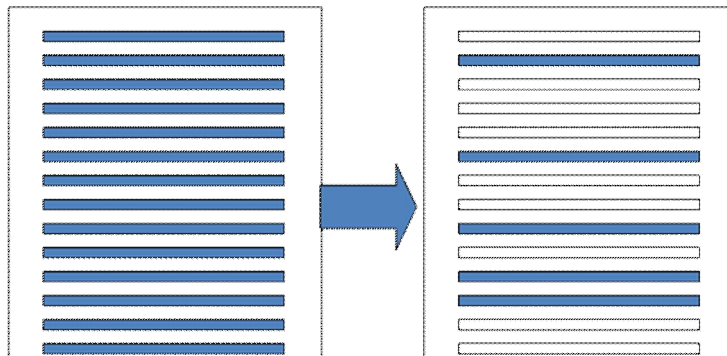
قبل از شرح بیشتر خلاصه‌سازی جایز دانستیم که ابتدا جایگاه خلاصه‌سازی را شرح دهیم و در خلال تعاریف اولیه و مشکلات خلاصه‌ساز را بیاوریم. اکنون بیش از 80 درصد از دانش ما به صورت متن، مستندات و دیگر صورت‌های رسانه نظیر ویدیو و صدا نگهداری می‌شود. اگر از منظر علوم کامپیوتری به این مستندات نگاه کنیم همه‌ی آن‌ها به طبعی غیر ساخت یافته وابسته‌اند. یک فرد برای دریافت دانش از اطلاعات یک متن، بایستی ابتدا آنرا درک کند و سپس آنرا پردازش نماید تا بفهمد چه معانی و مفاهیمی در آن موجود است، چه ارتباطی میان مفاهیم وجود دارد و از میان این مفاهیم کدام مهم و مفید است. چون قضاوت انسان بی‌ثبات است و تحت تأثیر عوامل محیطی، سلیقه فرد و دانش فرد مورد مطالعه می‌باشد، خلاصه‌های مختلفی از افراد مختلف می‌توان داشت. از مهمترین مشکلات مربوط به خلاصه‌سازی توسط انسان این است که هرچه تعداد مستندات بیشتر شود زمان

بیشتری را باید صرف مطالعه و خلاصه‌سازی کرد و از همه مهمتر سردرگمی کاربران در گزینش مطالب را می‌توان اشاره کرد.

## 2-3- خلاصه‌سازی خودکار

با گسترش روزافزون حجم اطلاعات موجود در وب و افزایش چشم‌گیر مقالات منتشر شده در زمینه‌های مختلف علمی، دسترسی درست و مطالعه مورد نیاز، همواره یکی از مشکلات محققان و پژوهشگران قرن 21 می‌باشد. اینکه چگونه با این حجم می‌توان در روز چندین کتاب را مطالعه نمود. آیا می‌توان سیستمی طراحی نمود که بتواند با داده‌های موجود به تمامی سوالات ما پاسخ دهد، اینها سوالاتی است که پاسخ آنها را می‌توان در یک سیستم خلاصه‌ساز متن جستجو کرد.

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]. بنابر تعریف ارائه شده در استاندارد ISO 215 [2] سال 1986، خلاصه "یک بازگویی مختصر از سند" می‌باشد (شکل 2-2).



شکل 2-2- خلاصه‌سازی خودکار سند

همانطور که اشاره شد خلاصه‌سازی خودکار توسط کامپیوتر انجام می‌شود و به همین دلیل تفاوت‌های زیادی با خلاصه‌سازی که توسط انسان تولید می‌شود دارد. انسان‌ها با توجه به هوش و شعور ذاتی خود قادر به درک و فهم مفاهیم موجود در متن و ارتباط بین آنها می‌باشند و این در حالی است که انجام این عملیات توسط ماشین کار بسیار دشوار و پیچیده‌ای می‌باشد. از طرفی دیگر، انسان‌ها با توجه به سطح دانش و پی‌زمینه‌ی اطلاعاتی که دارند دید متفاوتی از خلاصه‌ی یک متن یکسان دارند. به عنوان مثال کسی که سال‌ها در زمینه شبکه‌های کامپیوتری به تحقیق و مطالعه پرداخته است با کسی که به تازگی قصد تحقیق و مطالعه در زمینه‌ی شبکه‌های کامپیوتری را دارد، متفاوت بوده و خلاصه‌ای که این دو فرد از یک متن در زمینه‌ی شبکه‌های کامپیوتری تولید می‌کنند قطعاً یکسان نخواهد بود.

خلاصه‌سازی خودکار متن خود مشکلات و سختی‌های زیادی دارد اما مزیت‌های نسبت به خلاصه‌سازی دستی دارد که می‌توان به موارد زیر اشاره کرد:

- 1- ارتباط بین موجودیت‌های متن در خلاصه و موقعیت آنها در متن اصلی به آسانی قابل ایجاد است.
- 2- اندازه خلاصه قابل کنترل است.
- 3- محتوای آن قابل پیش‌بینی است.

این مطلب در چند سال اخیر بسیار مورد توجه واقع شده و منجر به ارائه بحث خلاصه‌سازی مبتنی بر کاربر یا خلاصه‌سازی شخصی‌سازی شده، گشته است. در ادامه به این مطلب بیشتر خواهیم پرداخت. سیستم‌های خلاصه‌ساز در دنیای امروز کاربردهای فراوانی دارند. تولید خلاصه‌های کتب مختلف و مقالات علمی، تولید خلاصه اخبار و انتقال آن از طریق سیستم‌های نظیر تلفن همراه، نمایش خلاصه سند یافته شده توسط موتور جستجو، تولید سیستم‌های پاسخ‌گویی به سوال و ... همگی از کاربردهای این سیستم می‌باشند.

## 2-4- فرایند خلاصه‌سازی

خلاصه‌سازی یکی از کاربردهای پردازش متن است. پردازش متن شامل چهار سطح است، پردازش لغوی، پردازش ساخت واژگی، پردازش نحوی و پردازش معنایی [1]. هر یک از این کاربردهای فراوان پردازش متن، از جمله بازیابی اطلاعات، خلاصه‌سازی، درک، تولید، ترجمه، پرسش و پاسخ، استخراج دانش از متون و موارد دیگر با توجه به گستردگی و پیچیدگی، در یک یا چند سطح فوق به انجام می‌رسد. خلاصه‌سازی، یکی از پیچیده‌ترین کاربردهای پردازش متن است و معمولاً با چند سطح از پردازش متن درگیر است. در ادامه به پردازش‌های لغوی و ساخت واژگی خواهیم پرداخت و پردازش‌های نحوی و پردازش‌های معنایی را به فراخور نیاز بعداً آنها را شرح خواهیم داد.

### 3-4-1- پردازش لغوی

منظور از پردازش لغوی شناسایی مرز لغات و جملات در یک متن است. این مرز ممکن است به شکل ساده توسط جداکننده‌هایی مانند فاصله، کاما، نقطه، علامت سؤال و ... تعیین شود و یا نیاز به پردازش‌های بیشتری داشته باشد، مانند زمانی که میان بخش‌هایی یک کلمه از فاصله استفاده می‌کنیم (مثل کلمه «می‌توان») و یا وقتی دو کلمه مجزا را بدون فاصله یا پی در پی می‌نویسیم (مثل عبارت در «برابر باد»). تعیین مرز کلمات در زبان فارسی به دلیل گوناگونی رسم الخط و عدم وجود استانداردهای نگارشی و همچنین به دلیل وجود شکل‌های مختلف حروف (اول - وسط - آخر چسبان و غیرچسبان) بیش از زبان انگلیسی مشکل‌ساز است. این مشکل در