

الحجامة



**دانشگاه صنعتی امیرکبیر (پلی تکنیک تهران)**

**دانشکده مهندسی کامپیوتر و فناوری اطلاعات**

پایان نامه کارشناسی ارشد

در رشته‌ی مهندسی کامپیوتر گرایش هوش ماشین و رباتیک

**طراحی و پیاده‌سازی یک سیستم استخراج اطلاعات**

**با استفاده از روش‌های تطبیقی**

**برای متون غیر ساخت‌یافته‌ی زبان فارسی**

نگارش:

**حمید خدایی**

استاد راهنما:

**دکتر محمدرضا مطش بروجردی**

بهمن ۱۳۸۷

باسپاس فراوان از جناب آقای دکتر  
بروجردی که یاری صمیمانه، بی دریغ، و از هر  
حیث‌شان انجام این پروژه را تسهیل نمود.

تقدیم بہ پدر و مادرم

### چکیده:

در پردازش زبان طبیعی، استخراج اطلاعات نوعی بازیابی اطلاعات بوده که قادر به استخراج اطلاعات ساختاریافته از متون فاقد ساختار به‌وسیله‌ی الگوهایی از پیش تعریف‌شده است. در این پروژه، یک روش شبه‌بی‌نظارت مبتنی بر خوشه‌بندی در دو مرحله برای تعیین مناسب‌بودن، و همچنین طبقه‌بندی الگوهای استخراج اطلاعات از متن براساس نوع آنها طراحی، و برای زبان فارسی پیاده‌سازی شده است. مرحله‌ی اول خوشه‌بندی بر اساس فرکانس تکرار الگوها و مرحله‌ی دوم بر طبق محتویات معنایی آنها انجام می‌شود که می‌تواند نوع اطلاعاتی را که هر الگو استخراج می‌کند نیز به صورت خودکار تعیین کند. آزمایش روش بر روی یک دامنه‌ی خاص (اخبار زلزله) و یک دامنه‌ی عمومی‌تر (اخبار سوانح) انجام یافته و در حالت دامنه‌ی خاص، چند روش دیگر نیز پیاده‌سازی شده، و نتایج آنها مورد مقایسه قرار گرفته است.

برای پردازش متن، از یک روش تکه‌سازی مبتنی بر ماشین بردار پشتیبان به همراه الگوریتمی برای تعیین محدوده‌ی جملات ساده و شکستن جملات مرکب به چند جمله‌ی ساده استفاده شده است که می‌تواند یک جمله را به گروه‌های دستوری تقسیم، و پس از آن، جملات مرکب را به جمله‌هایی ساده تبدیل نماید. همچنین روش تعیین جزء سخن براساس درخت تصمیم نیز استفاده شده است.

### کلمات کلیدی:

استخراج اطلاعات، یادگیری تطبیقی الگو، خوشه‌بندی، پردازش زبان طبیعی، تکه‌سازی متن، تعیین جزء سخن، زبان فارسی، متن بدون ساختار، تحلیل جملات مرکب

# فهرست مطالب

|           |   |    |
|-----------|---|----|
| ۱         | پیشگفتار.....                                   | ۱  |
| ۱-۱       | استخراج اطلاعات.....                            | ۱  |
| ۲-۱       | تفاوت استخراج و بازیابی اطلاعات.....            | ۳  |
| ۳-۱       | تاریخچه‌ی استخراج اطلاعات.....                  | ۳  |
| ۴-۱       | کاربردهای استخراج اطلاعات.....                  | ۶  |
| ۵-۱       | مسائل مطرح در سیستم‌های استخراج اطلاعات.....    | ۷  |
| ۱-۵-۱     | مسأله‌ی قابلیت انتقال.....                      | ۷  |
| ۲-۵-۱     | فرضیات هر روش.....                              | ۷  |
| ۳-۵-۱     | ورودی و خروجی در سیستم‌های استخراج اطلاعات..... | ۸  |
| ۴-۵-۱     | کشف الگوها.....                                 | ۹  |
| ۶-۱       | معماری سیستم‌های استخراج اطلاعات.....           | ۹  |
| ۷-۱       | روش‌های طراحی.....                              | ۹  |
| ۸-۱       | الگوریتم‌های استخراج اطلاعات.....               | ۱۰ |
| ۱-۱-۱     | یادگیری قوانین.....                             | ۱۲ |
| ۱-۱-۸-۱   | روش‌های بانظارت.....                            | ۱۲ |
| ۱-۱-۱-۸-۱ | یادگیری گزاره‌ای.....                           | ۱۲ |
| ۲-۱-۱-۸-۱ | یادگیری رابطه‌ای.....                           | ۱۳ |

- ۱۴.....۲-۱-۸-۱. روشهای شبه بی نظارت.
- ۱۷.....۲-۱-۱. یادگیری با استفاده از مدل های آماری.
- ۱۷.....۹-۱. سیستم های موجود.
- ۱۸.....۱-۹-۱. سیستم های مبتنی بر مهندسی دانش برای کشف الگوها.
- ۱۸.....۱-۱-۹-۱. FASTUS.
- ۱۸.....۲-۱-۹-۱. GE NLTOOLSET.
- ۱۹.....1-9-1-3. PLUM.
- ۱۹.....۴-۱-۹-۱. Proteus.
- ۱۹.....۲-۹-۱. سیستم های مبتنی بر یادگیری بانظارت برای کشف الگوها.
- ۱۹.....۱-۲-۹-۱. AutoSlog.
- ۱۹.....۲-۲-۹-۱. PALKA.
- ۲۰.....۳-۲-۹-۱. CRYSTAL.
- ۲۰.....۴-۲-۹-۱. LIEP.
- ۲۰.....۵-۲-۹-۱. WHISK.
- ۲۱.....۶-۲-۹-۱. RAPIER.
- ۲۱.....۷-۲-۹-۱. GATE.
- ۲۱.....۳-۹-۱. سیستم های مبتنی بر کشف شبه بی نظارت الگوها.
- ۲۱.....۱-۳-۹-۱. AutoSlog-TS.
- ۲۲.....۲-۳-۹-۱. Mutual Bootstrapping.
- ۲۳.....۳-۳-۹-۱. EXDISCO.
- ۲۳.....۴-۳-۹-۱. Snowball.

- ۲۴.....QDIE ۵-۳-۹-۱
- ۲۴..... ۱۰-۱. نحوه‌ی ارزیابی سیستم‌های استخراج اطلاعات
- ۲۶..... ۱۱-۱. پردازش زبان طبیعی
- ۲۶..... ۱-۱۱-۱. مقدمه
- ۲۷..... ۲-۱۱-۱. اجزای سیستم پردازش زبان طبیعی
- ۲۷..... ۱-۲-۱۱-۱. پیش‌پردازش
- ۲۷..... ۲-۲-۱۱-۱. تعیین جزء سخن
- ۳۰..... ۱-۲-۲-۱۱-۱. روش مدل مخفی مارکوف و الگوریتم ویتربی
- ۳۱..... ۲-۲-۲-۱۱-۱. روش درخت تصمیم
- ۳۲..... ۳-۲-۲-۱۱-۱. تعیین جزء سخن برای کلمات ناشناخته
- ۳۳..... ۳-۲-۱۱-۱. تعیین گروه‌های دستوری
- ۳۴..... ۱-۳-۲-۱۱-۱. ماشین بردار پشتیبان
- ۳۵..... ۲-۳-۲-۱۱-۱. تعیین ویژگی‌ها
- ۳۷..... ۴-۲-۱۱-۱. تحلیل نحوی و تفسیر معنایی
- ۳۸..... ۵-۲-۱۱-۱. تحلیل مباحثه‌ای
- ۳۸..... ۱۲-۱. خلاصه و جمع‌بندی
- ۴۰..... ۲. پیش‌پردازش
- ۴۰..... ۱-۲. پیش‌پردازش متن
- ۴۲..... ۲-۲. جداسازی جمله‌ها و واژه‌ها
- ۴۴..... ۳-۲. تعیین جزء سخن
- ۴۷..... ۴-۲. تعیین گروه‌های دستوری



|    |                                      |       |
|----|--------------------------------------|-------|
| ۵۰ | ..... تحلیل نحوی                     | ۵-۲   |
| ۵۱ | ..... استفاده از تحلیلهای نحوی       | ۱-۵-۲ |
| ۵۲ | ..... استفاده از تکه‌سازی در سطح دوم | ۲-۵-۲ |
| ۵۳ | ..... الگوریتم تحلیل نحوی            | ۳-۵-۲ |
| ۵۵ | ..... تحلیل جملات محذوف              | ۴-۵-۲ |
| ۵۶ | ..... تفسیر معنایی                   | ۶-۲   |
| ۵۸ | ..... خلاصه و جمع‌بندی               | ۷-۲   |
| ۶۰ | ..... یادگیری الگوها                 | ۳     |
| ۶۰ | ..... استخراج الگوها                 | ۱-۳   |
| ۶۲ | ..... فیلتر کردن الگوهای مناسب       | ۲-۳   |
| ۶۵ | ..... آزمایش و ارزیابی               | ۳-۳   |
| ۷۰ | ..... تغییر دامنه                    | ۴-۳   |
| ۷۱ | ..... خلاصه و جمع‌بندی               | ۵-۳   |
| ۷۲ | ..... نتیجه‌گیری                     | ۶-۳   |
| ۷۲ | ..... کارهای آینده                   | ۷-۳   |
| ۷۴ | ..... واژه‌نامه                      | ۴     |
| ۷۶ | ..... منابع و مراجع                  | ۵     |
| ۸۱ | ..... پیوست                          | ۶     |
| ۸۱ | ..... دستور زبان فارسی               | ۱-۶   |

- ۱-۱-۶ انواع نقش‌های دستوری..... ۱۲
- ۱-۱-۶-۱. فاعل..... ۸۲
- ۱-۱-۶-۲. مسندالیه..... ۸۲
- ۱-۱-۶-۳. مفعول..... ۸۲
- ۱-۱-۶-۴. مسند..... ۸۳
- ۱-۱-۶-۵. متمم..... ۸۳
- ۱-۱-۶-۶. قید..... ۸۴
- ۱-۱-۶-۷. فعل..... ۸۴
- ۱-۱-۶-۸. حروف عطف و ربط..... ۸۷
- ۱-۱-۶-۹. منادا..... ۸۸
- ۱-۱-۶-۱۰. حروف شرط..... ۸۸
- ۱-۱-۶-۱۱. صفت..... ۸۸
- ۱-۱-۶-۱۲. مضافالیه..... ۸۹
- ۲-۱-۶ انواع جمله..... ۱۹
- ۱-۲-۱-۶-۱. جمله‌ی اسنادی و فعلی..... ۸۹
- ۱-۲-۱-۶-۲. جمله‌ی ساده و جمله‌ی مرکب..... ۹۰
- ۱-۲-۲-۱-۶-۱. جمله‌ی ساده..... ۹۰
- ۱-۲-۲-۱-۶-۲. جمله‌ی مرکب..... ۹۰
- ۱-۲-۱-۶-۳. جمله‌ی کامل و محذوف..... ۹۲
- ۲-۶. گرامر زبان فارسی..... ۹۴

# فهرست جداول

---

- جدول ۱-۱: نمونه‌ای از اطلاعات استخراج شده ..... ۲
- جدول ۲-۱: دسته‌بندی روش‌های مبتنی بر یادگیری ماشین ..... ۱۱
- جدول ۳-۱: نمونه‌ای از اکتشاف‌های AUTOSLOG ..... ۱۹
- جدول ۴-۱: نتایج گزارش شده‌ی تعیین جزء سخن با روش درخت تصمیم برای زبان انگلیسی ..... ۳۲
- جدول ۵-۱: مثالی از IOB و IOE و نظایر آن‌ها ..... ۳۶
- جدول ۱-۲: نگاشت تبدیل کاراکترهای یونیکد به اسکی ..... ۴۱
- جدول ۲-۲: برچسب‌های پیکره بیجن خان ..... ۴۴
- جدول ۳-۲: نتایج آموزش تعیین‌کننده‌ی جزء سخن با روش درخت تصمیم برای زبان فارسی ..... ۴۶
- جدول ۴-۲: نتایج گزارش شده‌ی تعیین جزء سخن برای چند روش ..... ۴۷
- جدول ۵-۲: صحت الگوریتم تشخیص گروه‌های دستوری ..... ۵۰
- جدول ۶-۲: نتیجه‌ی تکه‌سازی در سطح دوم ..... ۵۳
- جدول ۷-۲: موجودیت‌های نامی تعریف شده ..... ۵۶
- جدول ۱-۳: درصد نمونه‌ها در هر خوشه در مرحله‌ی اول ..... ۶۷
- جدول ۲-۳: نتیجه‌ی تخصیص خوشه به شیارها ..... ۶۸
- جدول ۳-۳: نتیجه‌ی ارزیابی نهایی سیستم استخراج اطلاعات ..... ۶۸
- جدول ۴-۳: نتیجه‌ی ارزیابی نهایی روش AUTOSLOG-TS ..... ۶۹

جدول ۳-۵: نتیجه‌ی ارزیابی نهایی روش BASILI ..... ۶۹

جدول ۳-۶: نتیجه‌ی ارزیابی نهایی روش خوشه‌بندی با دامنه‌ی گسترده ..... ۷۰

# فهرست شکل‌ها

---

- شکل ۱-۱: معماری متداول یک سیستم استخراج اطلاعات (TURMO 2006) ..... ۹
- شکل ۱-۲: یک گره مفهوم که توسط AUTOSLOG استقرار می‌شود. (RILOFF 1996) ..... ۱۳
- شکل ۱-۳: فلوجارت AUTOSLOG-TS ..... ۲۲
- شکل ۱-۴: اجزای اصلی سیستم SNOWBALL ..... ۲۳
- شکل ۱-۵: مثال واژه‌نامه برای واژگان پسوندها (SCHMID 1994) ..... ۳۳
- شکل ۱-۲: معماری سیستم استخراج اطلاعات ..... ۴۰
- شکل ۲-۲: درخت لغات چند قسمتی ..... ۴۲
- شکل ۲-۳: درخت نحوی جمله‌ی «زلزله‌ای که ...» ..... ۵۵
- شکل ۱-۳: نتیجه‌ی خوشه‌بندی مرحله‌ی اول ..... ۶۷

## ۱. پیش‌گفتار

### ۱-۱. استخراج اطلاعات

داده‌های فاقد ساختار، مانند متن ساده، معمولاً بخش اعظم دانش ذخیره‌شده‌ی یک سازمان را تشکیل می‌دهد. اما دسترسی، تحلیل، جستجو، و یا استفاده از آن‌ها توسط افراد، بسیار پرهزینه و زمان‌بر می‌باشد. بنابراین، چنانچه ابزارهایی موجود باشد که بتوان به‌وسیله‌ی آن‌ها، اطلاعات موجود در متن‌ها را به صورت یک ساختار ساده و قابل فهم توسط ماشین استخراج کرد، می‌توان از آن اطلاعات به نحو مطلوب استفاده نمود. روشی که برای این منظور به کار گرفته می‌شود، متن‌کاوی<sup>۱</sup> نام دارد.

متن‌کاوی که زیر شاخه‌ای از داده‌کاوی<sup>۲</sup> است، دانشی است درباره‌ی نحوه‌ی جستجوی الگوها در متن به زبان طبیعی که به صورت «کشف اطلاعات به‌وسیله‌ی کامپیوتر و با استفاده از استخراج اطلاعات از منابع نوشتاری مختلف [Hearst 2004]» تعریف شده است.

استخراج اطلاعات<sup>۳</sup> یکی از برجسته‌ترین تکنیک‌هایی است که در متن‌کاوی مورد استفاده قرار می‌گیرد. به طور خاص، در این فن‌آوری می‌توان با استفاده از ترکیب روش‌ها و تکنیک‌های پردازش زبان طبیعی<sup>۴</sup> کارایی بالایی را برای کاوش متن از دامنه‌های<sup>۵</sup> گوناگون به دست آورد. بنابراین، استخراج اطلاعات، یک فن‌آوری منشعب از پردازش زبان طبیعی است که به تحلیل متن فاقد ساختار و به زبان طبیعی برای مشخص نمودن اطلاعات یا وقایعی می‌پردازد که به صورت صریح یا ضمنی در آن وجود دارد [Cowie 1996].

یکی از دلایل تمایل پژوهشگران به این دانش، نقش آن در فراهم نمودن امکان مقایسه و ارزیابی تکنیک‌های مختلف پردازش زبان طبیعی می‌باشد [Cowie 1996]. علاوه بر آن، این فناوری،

---

<sup>1</sup> Text Mining (TM)

<sup>2</sup> Data Mining

<sup>3</sup> Information Extraction (IE)

<sup>4</sup> Natural Language Processing (NLP), Natural Language Understanding (NLU)

<sup>5</sup> Domain

می‌تواند کاربردهای گوناگونی برای صنایع مختلف به خصوص فعالان بورس، بانک‌ها، ناشران، و حتی دولتمردان داشته باشد. به طور کل، هر شخص یا سازمانی که نیاز به مطالعه‌ی روزانه‌ی اخبار نوشتاری برای استخراج اطلاعات از آن داشته باشد، می‌تواند از این گونه سیستم‌ها استفاده نماید.

یکی دیگر از تعاریفی که می‌توان برای استخراج اطلاعات در نظر گرفت این است که در این دانش سعی بر آن است که بتوان اطلاعات موجود در یک متن به زبان طبیعی را به شکل ساختاری سازمان‌یافته و جدول‌گونه<sup>۱</sup> ارائه داد [Grishman 1997]. منظور از ساختار جدولی این است که پیش از اعمال یک الگوریتم استخراج اطلاعات، جدولی با دو ستون موجود باشد که در ستون اول، نوع اطلاعاتی که انتظار می‌رود الگوریتم آنها را بیابد قرار دارد و ستون دوم، پس از اجرای الگوریتم، با اطلاعات استخراج‌شده پر گردد. به هر ردیف، که شامل یک زوج نوع اطلاعات و اطلاعات استخراج‌شده است یک شیار<sup>۲</sup> یا شکاف گفته می‌شود. برای نمونه، اگر خبری مانند خبر زیر موجود باشد، سیستم استخراج اطلاعات باید بتواند اطلاعات را به صورت جدول ۱-۱ استخراج نماید:

### خبر: انفجار مرکز نیویورک را به لرزه درآورد<sup>۳</sup>

انفجار شدیدی در ساعات پیر رفت‌وآمد بعد از ظهر در مرکز نیویورک در اثر ترکیدن یک لوله تاسیسات در دل زمین باعث مرگ یک نفر و زخمی شدن دست کم ۱۸ نفر شده است. پلیس تعداد زیادی از ساختمان‌های اطراف و ایستگاه مرکزی شهر را که انفجار در نزدیکی آن روی داد تخلیه کرد.

انفجار در اثر ترکیدن یک لوله بخار در زیر خیابان لکزینگتون در ناحیه منهتن روی داد. اداره پلیس نیویورک گفت حادثه که باعث ایجاد حفره بزرگی در خیابان شد، تروریستی نبوده است.

انفجار بزرگ روز چهارشنبه که درست قبل از ساعت شش بعد از ظهر روی داد باعث ایجاد وحشت و هرج و مرج در خیابان‌های بخش شرقی منهتن شد و درحالی که بخار با غلظت زیاد از دل زمین بیرون می‌زد مردم از صحنه می‌گریختند. ماموران آتش‌نشانی و سازمان‌های اضطراری به نقطه حادثه شتافته و بخش‌هایی از خیابان، در فاصله ایستگاه "گراند سنترال" و ساختمان کرایسلر را بستند. هزاران عابر از ایستگاه تخلیه شدند.

جدول ۱-۱: نمونه‌ای از اطلاعات استخراج‌شده

| انفجار                              | نوع حادثه |
|-------------------------------------|-----------|
| ترکیدن یک لوله تاسیسات در دل زمین   | علت       |
| مرگ یک نفر و زخمی شدن دست کم ۱۸ نفر | تلفات     |
| خیابان لکزینگتون، منهتن، نیویورک    | مکان      |
| چهارشنبه-قبل از ساعت ۶ بعد از ظهر   | زمان      |

<sup>۱</sup> Tabular

<sup>۲</sup> Slot

<sup>۳</sup> منبع خبر BBC است.

به طوری که در مثال فوق مشاهده می‌شود، به‌ازای یک خبر، ۵ شیار اطلاعاتی نوع، زمان، مکان، تلفات، و علت حادثه استخراج گردیده است.

## ۲-۱. تفاوت استخراج و بازیابی اطلاعات

به طور سنتی، اطلاعات موجود در متون توسط افراد خبره‌ی مربوط به دامنه‌ی متن و به صورت دستی<sup>۱</sup> استخراج می‌شده‌است که بدیهی‌ست این عمل بسیار پرهزینه چه از لحاظ زمانی و چه از لحاظ مسائل مربوط به نیروی انسانی می‌باشد. به همین دلیل، در دهه‌های اخیر، تلاش‌هایی برای ایجاد سیستم‌هایی موسوم به هوش متنی<sup>۲</sup> انجام گرفته که هدف آن‌ها تغییر متن بوده چنان که بتوان به اطلاعات موجود در آن، دست یافت [JACOBS 1992, Turmo 2006]. متن‌ها معمولاً هنگامی که برای استفاده‌ی ماشینی تولید می‌شوند (مانند زبان‌های برنامه‌نویسی)، به صورت کاملاً ساختاریافته هستند ولی هنگامی که برای استفاده توسط انسان تهیه می‌شوند، به شکل زبان طبیعی بوده که فاقد ساختاری صریح می‌باشد. که این امر موجب می‌شود استخراج اطلاعات از آن‌ها نیاز به دانشی وسیع از زبان و دامنه‌ی متن داشته باشد.

هوش متنی به دو زیرشاخه‌ی اصلی بازیابی اطلاعات و استخراج اطلاعات، تقسیم شده است. در بازیابی اطلاعات، هدف این است که از مجموعه‌ای از مستندات، آنهایی که مرتبط با یک جستار<sup>۳</sup> ورودی‌اند، انتخاب شوند. در حالی که در استخراج اطلاعات، باید از داخل یک متن، اطلاعات مهم آن استخراج شده و به کاربر ارائه شود.

در سیستم‌های بازیابی اطلاعات، هر متن، تنها به صورت مجموعه‌ای از کلید-واژه‌ها پنداشته می‌شود. در حالی که در استخراج اطلاعات، به ساختار زبانی جملات توجه شده و اطلاعات بر اساس الگوهایی مطابق با ساختارهای دستوری زبان طبیعی متن استخراج می‌شوند.

در سیستم‌های بازیابی اطلاعات، برخلاف استخراج اطلاعات، از دانش پردازش زبان طبیعی، تنها برای مسائلی نظیر تعیین محدوده‌ی واژه‌ها، تعیین ریشه‌ی کلمه‌ها، تشکیل گروه‌های اسمی، و نظایر آن‌ها استفاده می‌شود. ولی در استخراج اطلاعات، سطوح بالاتر این دانش مانند تحلیل‌گر<sup>۴</sup>های نحوی، تفسیرکننده‌های معنایی نیز به کار گرفته می‌شوند.

## ۳-۱. تاریخچه‌ی استخراج اطلاعات

توسعه‌ی دانش استخراج اطلاعات وابسته به کنفرانس درک متن<sup>۵</sup> (MUC) بوده که از سال ۱۹۸۷ تا ۱۹۹۸ برگزار می‌شده‌است. هدف اصلی کلی کنفرانس مذکور، برگزاری رقابتی بر اساس

<sup>۱</sup> Manually

<sup>۲</sup> Text-based Intelligence (TBI)

<sup>۳</sup> Query

<sup>۴</sup> Parser

<sup>۵</sup> Message Understanding Conference (MUC): [http://www.itl.nist.gov/iaui/894.02/related\\_projects/muc/](http://www.itl.nist.gov/iaui/894.02/related_projects/muc/)



سیستم‌های استخراج اطلاعات از متون با دامنه‌ی محدود و قالب آزاد بوده است که توسط تیم‌های مختلف طراحی شده بودند. معمولاً برای هر کنفرانس، یک دامنه انتخاب می‌شد. خلاصه‌ای از فعالیت‌های انجام‌شده در دوره‌های مختلف، بدین شرح است:

- **MUC-1 (۱۹۸۷):** اولین MUC به صورت تحقیقاتی برگزار شد. به این صورت که نه الگوریتمی ارائه گردید و نه معیاری برای ارزیابی؛ و هر گروه، به طراحی قالب خود برای ارائه‌ی اطلاعات بسنده نمود. دامنه‌ی متون، عملیات تاکتیکی دریایی انتخاب شده بود.
- **MUC-2 (۱۹۸۹):** در این کنفرانس، دامنه‌ی کنفرانس قبلی انتخاب شد. الگوریتمی نیز برای پر کردن الگوها ارائه گردید. ۱۰ شیار نوع رویداد، عامل، زمان و مکان، اثر، و .. برای پر کردن به شرکت‌کنندگان معرفی شد. برای ارزیابی و شیوه‌ی آن، معیار واحدی ارائه نگردید.
- **MUC-3 (۱۹۹۱):** دامنه‌ی مستندات به وقایع تروریستی در آمریکای لاتین با ۱۸ شیار نوع حادثه، تاریخ، زمان، عامل، هدف، وسیله، و ... تغییر کرد. تعداد ۱۳۰۰ متن آموزشی و ۳۰۰ متن آزمایشی در اختیار شرکت‌کنندگان قرار گرفته بود. ۴ میزان<sup>۱</sup> برای اندازه‌گیری معرفی شده بود: تعداد شکاف‌هایی که درست استخراج شده‌اند (COR)، تعداد شکاف‌هایی که خالی مانده‌اند (INC)، تعداد شکاف‌هایی که اشتباه پر شده‌اند (SPUR)، تعداد اطلاعات موجود در متن که استخراج نشده‌اند (MISS)، و تعداد اطلاعاتی که ناقص استخراج شده‌اند (PAR). همچنین دو معیار بازخوانی (R)<sup>۲</sup> و دقت (P)<sup>۳</sup> برای ارزیابی به صورت زیر تعریف شده بودند:

$$R = \frac{COR + 0.5 \times PAR}{COR + PAR + INC + MISS} \quad (1)$$

$$P = \frac{COR + 0.5 \times PAR}{COR + PAR + INC + SPUR} \quad (2)$$

- **MUC-4 (۱۹۹۲):** این کنفرانس مشابه MUC-3 بود با این تفاوت که تعداد شکاف‌ها به ۲۴ افزایش یافت. برای ارزیابی نیز معیار F بر اساس دقت و بازخوانی به صورت زیر تعریف گردید:

$$F = \frac{(\beta^2 + 1) \times P \times R}{\beta^2 \times P + R} \quad 0 < \beta \leq 1 \quad (3)$$

- **MUC-5 (۱۹۹۳):** در کنفرانس‌های قبلی، هدف تنها استخراج اطلاعات از متون به زبان انگلیسی بود. در MUC-5، برای زبان ژاپنی نیز سیستمی ارائه شد. همچنین دو دامنه‌ی مختلف اخبار مالی در زمینه‌ی سرمایه‌گذاری مشترک و پیشرفت‌های محصولات میکرو-الکترونیک، پیشنهاد شد. برای ارزیابی، علاوه بر روش‌های کنفرانس‌های قبلی، معیارهایی نیز براساس خطا معرفی شدند.

<sup>1</sup> Measure

<sup>2</sup> Recall (R)

<sup>3</sup> Precision (P)

- **MUC-6 (۱۹۹۵):** در این کنفرانس، دامنه مالی مورد استفاده قرار گرفت و اهدافی نظیر استقلال از دامنه، قابلیت انتقال<sup>۱</sup>، و فهم عمیق تر متن عنوان شد. و دو معیار بازخوانی و دقت به صورت زیر تغییر کردند:

$$R = \frac{COR}{COR + INC + MISS} \quad (۴)$$

$$P = \frac{COR}{COR + INC + SPUR} \quad (۵)$$

- **MUC-7 (۱۹۹۸):** در این کنفرانس، دامنه‌ی سوانح هوایی استفاده شد و زبان چینی نیز مورد بررسی قرار گرفت.

به موازات MUC، اتحادیه‌ی اروپا بر روی برنامه‌ی LRE<sup>۲</sup> که در مورد روش‌ها و ابزارهای استخراج و بازیابی اطلاعات است، سرمایه‌گذاری نمود که نتایج آن، پروژه‌های ECRAN<sup>۳</sup>، SPARKLE<sup>۴</sup>، FACILIE<sup>۵</sup> و AVENTINUS<sup>۶</sup> می‌باشد.

پس از کنفرانس‌های MUC، از سال ۱۹۹۹، پژوهش بر روی فناوری استخراج اطلاعات به برنامه‌ی TIDES<sup>۷</sup> با سرمایه‌گذاری DARPA منتقل شد. تعدادی از پروژه‌های نهایی شامل RIPTIDES<sup>۸</sup>، PROTEUS<sup>۹</sup>، CREST<sup>۱۰</sup>، Coreference.com<sup>۱۱</sup> و UMass<sup>۱۲</sup> می‌باشد.

رقابت اولیه‌ی ACE، TIDES<sup>۱۳</sup> نام داشت که از سال ۱۹۹۹ با هدف استخراج خودکار محتوا از متون به زبان طبیعی، شروع به کار کرد. منابع استفاده‌شده شامل مطالبی بود که از منابعی نظیر مجلات الکترونیکی، اخبار همگانی، و روزنامه‌ها (از طریق OCR) گردآوری شد بود. فعالیت‌های انجام‌شده از طریق اضافه کردن موجودیت‌های سلسله‌مراتبی پیچیده‌تر از MUC شد.

در فاز دوم ACE (۲۰۰۱ تا ۲۰۰۲)، ارزیابی بر میزان تشخیص روابط و توصیفات اضافه شد. رقابت ۲۰۰۳، برای اولین بار شامل زبان‌های غیر انگلیسی (چینی و عربی) نیز شد. در سال ۲۰۰۴، برنامه به توضیح روابط عربی و نرمال‌سازی زمان و تاریخ گسترش یافت. همچنین یک الگوریتم برای تشخیص رویدادها ارائه شد. مفهوم رویداد در ACE ساده‌تر از MUC است و شامل موجودیت‌ها،

<sup>1</sup> Portability

<sup>2</sup> Linguistic Research and Engineering (<http://www.echi.lu/language/en/lehome.html>)

<sup>3</sup> <http://www.dcs.shef.ac.uk/intranet/research/networks/Ecran/>

<sup>4</sup> <http://www.informatics.susx.ac.uk/research/nlp/sparkle/sparkle.html>

<sup>5</sup> <http://tcc.itc.it/research/textec/projects/facile/facile.html>

<sup>6</sup> <http://www.dcs.shef.ac.uk/nlp/funded/aventinus.html>

<sup>7</sup> Tran-lingual Information Detection, Extraction, and Summarization: <http://www.darpa.mil/ipto/programs/tides/>

<sup>8</sup> <http://www.cs.cornell.edu/home/cardie/tides>

<sup>9</sup> <http://nlp.cs.nyu.edu/>

<sup>10</sup> <http://crl.nmsu.edu/Research/Projects/Crest>

<sup>11</sup> <http://www.coreference.com/lingpipe/>

<sup>12</sup> <http://ciir.cs.umass.edu/research/tides.html>

<sup>13</sup> Automatic Content Extraction: <http://www.nist.gov/speech/tests/ace/>

مقادیر، و زمان می‌باشد. در سال ۲۰۰۵، منابع جدیدی به هر سه زبان اضافه شد (مکالمات تلفنی، گروه‌های خبری، وبلاگ‌ها و ...).

علاوه بر کنفرانس‌های فوق، تلاش‌های متفرقه‌ی دیگری هم در زمینه‌ی دامنه‌های دیگر انجام شده است. به عنوان مثال، [Soderland 1995] سیستمی را برای استخراج اطلاعات پزشکی از گزارش‌های تریخیص بیمارستان‌ها طراحی نمود و [Holowczak 1997] از اطلاعات استخراج‌شده‌ی متون حقوقی، برای دسته‌بندی آنها استفاده کرد. همچنین [Glasgow 1998] برای دامنه‌ی بیمه‌ی عمر سیستمی را ارائه نمود.

برخی از تلاش‌ها نیز بر روی کاربردهای استخراج اطلاعات بر روی متون نیمه‌ساخت‌یافته انجام گرفت. مانند دامنه‌ی آگهی اجاره‌ی آپارتمان توسط [Soderland 1999]، دامنه‌ی biomedical توسط [Craven 1999].

به موازات TIDES، اتحادیه‌ی اروپا بر روی 'Pascal Network of Excellence' سرمایه‌گذاری نمود. این سازمان همراه با پروژه‌ی 'Dot.Kom European Challenge on Evaluation of Machine Learning for Information Extraction from Documents' را اجرا می‌کند. اولین رقابت آن در ژوئن ۲۰۰۴ شروع و ارزیابی‌های رسمی در نوامبر ۲۰۰۴ انجام گرفت. تفاوت عمده‌ی آن با رقابت‌های MUC در مبتنی بر یادگیری ماشین بودن آن است.

## ۴-۱. کاربردهای استخراج اطلاعات

از آنجایی که استخراج اطلاعات قادر است یک متن را به ساختار جدولی تبدیل کند، می‌توان از نتیجه‌ی آن، در هرگونه عملیات پردازشی استفاده نمود. زیرا جدول، ساختاری است که به سادگی قابل پردازش می‌باشد. به عنوان مثال، می‌توان اطلاعات بازار بورس را در یک دوره‌ی مشخص از روی اخبار موجود استخراج کرد و پردازش‌های لازم را بر روی آن انجام داد. پایگاه‌های داده‌ی رابطه‌ای براساس ساختاری جدولی تعریف می‌شوند. از این رو به‌وسیله‌ی یک سیستم استخراج اطلاعات، تبدیل متون به یک پایگاه داده‌ی حاوی اطلاعات مورد نیاز امکان‌پذیر است.

در کاربرد خلاصه‌سازی، باید بخش‌هایی از متن که حاوی اطلاعات مهم است مشخص شود. از این رو می‌توان سیستم استخراج اطلاعاتی را به کار گرفت که جملات دارای اطلاعات را برچسب‌گذاری کند. زیرا معمولاً جملاتی مهم هستند که دارای اطلاعات باشند.

در زمینه‌ی بازیابی اطلاعات نیز استخراج اطلاعات مورد استفاده قرار گرفته است. زیرا با استفاده از این دانش می‌توان به بخش‌هایی از متن بچسب‌ها معنایی تخصیص داد که باعث می‌شود

<sup>1</sup> <http://www.pascal-network.org/>

<sup>2</sup> <http://www.dot-kom.org>

استفاده از جستارهای غنی تر میسر گردد. چرا که امکان ارائه‌ی آن را به زبان طبیعی میسر می‌سازد. به عنوان مثال با دادن جستار «چه تعداد زلزله‌ی بالای ۶ ریشتر در دو سال گذشته در اندونزی رخ داده است» می‌توان در مستندات مربوط به زلزله که قبلاً اطلاعات آنها شامل «شدت، زمان، مکان، و ...» استخراج شده است و پس از فهم کلمه‌ی «چه تعداد»، اطلاعات مورد نیاز را با یک جستجوی بسیار ساده به کاربر ارائه داد.<sup>۱</sup>

در سیستم‌های پاسخ به پرسش که وظیفه دارند با خواندن یک متن به پرسش‌های کاربر در مورد آن پاسخ گویند نیز می‌توان از استخراج اطلاعات استفاده کرد. زیرا اغلب پرسش‌های کاربر در مورد اطلاعات موجود در متن است و اگر این اطلاعات از قبل استخراج شده باشند، پاسخ به پرسش‌های کاربر به‌سادگی میسر خواهد بود.

در سیستم‌های اخذ دانش<sup>۲</sup> از متن، می‌توان از استخراج اطلاعات به‌منظور فراگیری دانش از نوع واقعیت‌ها<sup>۳</sup> استفاده نمود. زیرا هر شیار استخراج شده مانند  $\langle \text{Slot}, \text{Value} \rangle$  می‌تواند به صورت منطقی  $\text{Slot}(\text{Value})$  در نظر گرفته شود.

## ۵-۱. مسائل مطرح در سیستم‌های استخراج اطلاعات

### ۱-۵-۱. مساله‌ی قابلیت انتقال<sup>۴</sup>

یکی از مسائلی که در کاربردی‌شدن استخراج اطلاعات نقش فراوانی دارد، مساله‌ی قابلیت انتقال یا هزینه‌ی تطبیق دادن یک سیستم با دامنه‌های جدید می‌باشد. از آنجایی که سیستم استخراج اطلاعات شدیداً وابسته به دامنه می‌باشد، لذا برای انتقال به یک دامنه‌ی جدید می‌بایست سیستم جدیدی ایجاد نمود که این عمل، دارای هزینه است و یکی از عوامل موثر بر این هزینه، میزان تغییراتی است که باید در سیستم داده شود. عامل دیگر، مقدار اطلاعاتی است که باید مجدداً از کاربر دریافت گردد. هر چه این هزینه‌ها کمتر باشد، سیستم از این لحاظ مطلوب‌تر خواهد بود.

### ۲-۵-۱. فرضیات هر روش

برای انتخاب یا طراحی یک الگوریتم استخراج اطلاعات و این که بهتر است از روش‌های دستی استفاده شود و یا روش‌های مبتنی بر یادگیری ماشین، فرضیات زیر باید مد نظر قرار گیرند:

۱. موجود بودن داده‌های آموزشی: اگر داده‌های آموزشی در دسترس یا دستیابی به آن کم‌هزینه باشد، روش‌های یادگیری می‌توانند مورد استفاده قرار گیرند.

<sup>۱</sup> ر.ک. [Moens 2006]

<sup>۲</sup> Knowledge Acquisition

<sup>۳</sup> Fact

<sup>۴</sup> Portability