

# دانشگاه کامپیوتن

## dataacademy.ir

بکارگیری تکنیکی از داده‌کاوی برای بهبود کارآیی  
سیستم‌های استخراج اطلاعات از متن

مجتبی شکری<sup>۲</sup>

احمد عبداللهزاده<sup>۱</sup>

۱ تهران - خ سمهیه - دانشگاه صنعتی امیرکبیر - دانشکده کامپیوتر و فناوری اطلاعات

۲ تهران - خ کارگر شمالی - ساختمان شماره ۲۴۰ - طبقه پنجم - شرکت راهبر

ahmad@ce.aut.ac.ir

### چکیده:

با رشد سریع حجم متون قابل دسترس بهخصوص در شبکه‌ی جهانی اینترنت کار استخراج اطلاعات از متن اهمیت ویژه‌ای یافته است. خروجی سیستم استخراج اطلاعات، پایگاه داده‌ای از اطلاعات موجود در متن می‌باشد. از طرفی تکنیک‌های داده‌کاوی این امکان را به ما می‌دهند که دانشی را به شکل مجموعه‌ای از قواعد از محتویات پایگاه داده و ارتباط بین آن‌ها بدست آوریم. در این مقاله از تکنیک تولید قواعد تداعی گر در داده‌کاوی استفاده می‌کیم. بدین منظور با استفاده از الگوریتم استاندارد C4.5RULES و اعمال آن بر روی خروجی سیستم استخراج اطلاعات، دانشی (مجموعه‌ی قواعد پیشگو) بدست می‌آوریم که از آن برای بالاتر بردن احتمال درستی اطلاعات استخراج شده، پیشگویی اطلاعات استخراج نشده و مقایسه بین قواعد استخراج اطلاعات استفاده می‌کیم. در ادامه، طرحی برای معماری سیستم‌های استخراج اطلاعات پیشنهاد می‌کنیم تا بتواند از این تکنیک برای بهبود کارآیی‌شان استفاده کنند. سپس نتایج پیاده‌سازی این تکنیک را بر روی سیستم WHISK (یکی از سیستم‌های باز استخراج اطلاعات موجود) مورد ارزیابی قرار می‌دهیم.

**واژه‌های کلیدی:** استخراج اطلاعات از متن - داده‌کاوی - وب‌کاوی - استخراج دانش از پایگاه داده - معماری سیستم‌های استخراج اطلاعات از متن

### ۱- مقدمه

زمانیکه بخواهیم اطلاعات خاصی را از انبوه متون در دسترس استخراج کنیم، اهمیت سیستم‌هایی که بتوانند متون را به صورت خودکار پردازش کنند، بیشتر مشخص خواهد شد. در وهله اول باید بتوانیم از انبوه مستندات موجود، مستنداتی را انتخاب کنیم که مربوط به زمینه‌ی مورد نظر ما باشد. این وظیفه به عهده‌ی سیستم‌های بازیابی اطلاعات است. در مرحله‌ی دوم باید در این مستندات جستجو کرده و اطلاعات مورد نیازمان را جمع‌آوری کنیم. اینجاست که سیستم‌های استخراج اطلاعات به کمکمان می‌آیند. استخراج اطلاعات از متن شامل ارائه‌ی یک قالب ساختارمند (مانند یک پایگاه داده) از اطلاعات دلخواه موجود در متن می‌باشد [۱]. دو معیار مهم برای ارزیابی کارآیی سیستم‌های استخراج اطلاعات وجود دارد. اول اینکه، چه درصدی از اطلاعات استخراجی صحیح هستند و دوم اینکه، چه درصدی از اطلاعات موجود در متن استخراج یافته‌اند [۲].

<sup>۱</sup> دانشیار، هیئت علمی دانشکده کامپیوتر دانشگاه صنعتی امیرکبیر

<sup>۲</sup> کارشناسی ارشد نرم افزار دانشگاه صنعتی امیرکبیر، هیئت علمی دانشگاه آزاد واحد دماوند

# فناوری اطلاعات و دانش

تهران / دانشگاه صنعتی امیرکبیر / ۳-۵ خرداد ۱۳۸۴

بکارگیری تکنیک‌های دیگری که در هوش مصنوعی وجود دارند، می‌تواند منجر به بهبود معیارهای کارآیی سیستم‌های استخراج اطلاعات شود [۳]. موضوع مقاله مزبور استفاده از تکنیک تولید قواعد تداعی‌گر در داده‌کاوی برای کاراتر کردن سیستم‌های استخراج اطلاعات است. در داده‌کاوی از روش‌های آماری یا یادگیری ماشین برای کشف روابط تازه در پایگاه‌های داده‌ی رابطه‌ای بهره برده می‌شود. در سیستم استخراج اطلاعات از مستندات زبان طبیعی، اطلاعات دلخواه استخراج می‌شود. به عبارت دیگر از متن بدون ساختار، اطلاعات ساختارمند مانند یک پایگاه داده استخراج می‌شود. پیوند داده‌کاوی و سیستم استخراج اطلاعات زمانی است که خروجی سیستم استخراج اطلاعات را به عنوان ورودی برای داده‌کاوی در نظر بگیریم و با استفاده از تکنیک‌های موجود در داده‌کاوی، روابط تازه‌ای در اطلاعات استخراجی کشف کنیم. حاصل ادغام استخراج اطلاعات از متن و داده‌کاوی، استخراج دانش از متن خواهد بود. حال دانشی در اختیار ماست که هم خود یک خروجی بالارزش است و هم می‌توان از آن برای منظورهای دیگر مانند بالا بردن کارآیی سیستم‌های استخراج اطلاعات استفاده کرد. ما از تکنیک‌های موجود در داده‌کاوی از روش C4.5RULES [۴] برای تولید قواعد تداعی‌گر استفاده خواهیم کرد که یک روش استقرای قاعده با استفاده از ساخت درخت تصمیم از روی محتویات پایگاه‌های داده می‌باشد. سیستم WHISK [۱] که یکی از سیستم باز استخراج اطلاعات از متن می‌باشد را برای اضافه کردن تکنیک مورد نظر انتخاب کرده‌ایم و نتایج حاصله را مورد ارزیابی و تجزیه و تحلیل قرار داده‌ایم.

## ۲- پیش‌زمینه

در این بخش در ابتدا درباره سیستم‌های استخراج اطلاعات بطور کلی و سیستم WHISK بطور خاص صحبت می‌کنیم و سپس درباره داده‌کاوی و روش استقرای قاعده C4.5RULES صحبت خواهیم کرد.

## ۱- سیستم‌های استخراج اطلاعات

وظیفه سیستم استخراج اطلاعات، ارائه اطلاعات دلخواه، بصورت ساختارمند از اطلاعات موجود در متن است. نوع متن ورودی به سیستم‌های استخراج اطلاعات در سه دسته‌ی ساختارمند، شبیه ساختارمند و بدون ساختار تقسیم‌بندی می‌شوند [۵]. متن ساختارمند از یک قالب خاص پیروی می‌کند و از یک واژگان محدود استفاده می‌کند مانند گزارشات آزمایشگاهی. متن شبیه ساختارمند دارای قالب کاملاً مشخصی نمی‌باشد ولی مانند متن بدون ساختار نیز نیست و اغلب بصورت اطلاعات غیر گرامی و تلگرافی می‌باشد. متن بدون ساختار شامل جملات زبان طبیعی مانند داستان‌ها می‌باشد. سیستم استخراج اطلاعات برای استخراج اطلاعات از متن احتیاج به دانشی دارد که اغلب به صورت مجموعه قواعد استخراج اطلاعات می‌باشد که الگویی را در متن جستجو می‌کند و خروجی را بر اساس آن بدست می‌آورد [۶]. در سال‌های اخیر تمرکز تحقیقات بر روی استفاده از روش‌های آماری و یادگیری ماشین برای بدست آوردن مجموعه قواعد استخراج اطلاعات به صورت خودکار بوده است. سیستم‌های استخراج اطلاعات را می‌توان به عنوان پس‌پردازش بر روی خروجی بدست آمده از سیستم‌های بازیابی اطلاعات (مانند موتورهای جستجو) استفاده کرد و اطلاعات دقیق‌تری بدست آورد [۷].

زمانی یک سیستم استخراج اطلاعات را موفق می‌نامیم که به این اهداف رسیده باشد [۷]:

- ۱- جملاتی از متن را باید که حاوی اطلاعات مرتبط باشد.
- ۲- اطلاعات مرتبط با زمینه مورد علاقه را استخراج کند.
- ۳- اطلاعات را به هم پیوند دهد و در یک قالب از پیش تعیین شده در خروجی ارائه کند.
- ۴- از اطلاعات غیر مرتبط صرف نظر کند (بسیار مهم).

## ۲- سیستم WHISK

در سال‌های اخیر چندین سیستم استخراج اطلاعات ساخته شده است که اغلب از روش‌های یادگیری ماشین بهره برده‌اند. از جمله مهمترین آنها می‌توان به سیستم‌های AutoSlog در سال ۹۳، CRYSTAL در سال ۹۵ WIEN در سال ۹۵

# فناوری اطلاعات و دانش

تهران / دانشگاه صنعتی امیرکبیر / ۳-۵ خرداد ۱۳۸۴

RAPIER در سال ۹۷ SRV در سال ۹۸ WHISK در سال ۹۸ و STALKER در سال ۹۹ اشاره کرد.<sup>[۸]</sup>

در سیستم WHISK برای بدست آوردن مجموعه قواعد استخراج اطلاعات از یادگیری با نظارت استفاده شده است. قواعد استخراج اطلاعات در WHISK بر پایه‌ی عبارات باقاعدۀ می‌باشد که در آن الگویی ارائه می‌شود که باید در متن جستجو شده و در صورت تطابق آن با یک قسمت از متن، خروجی لازم بدست می‌آید. برای نمونه در شکل ۱ باقاعدۀ برای استخراج اطلاعات از آگهی‌های اجاره خانه آمده است. الگویی را نشان می‌دهد که باید در متن جستجو شود و Output، خروجی مربوطه را نشان می‌دهد.<sup>[۱]</sup>

**Pattern:: \* (Nghbr) \* (Digit) ‘Bdrm’ \* ‘\$’ (Number)**

**Output:: Rental {Neighborhood \$1} {Bedrooms \$2} {Price \$3}**

شکل ۱: نمونه‌ای از باقاعدۀ استخراج اطلاعات از متن در WHISK

پارامترهای ارائه شده در خروجی از عبارات مشخص شده در الگو بدست می‌آید. بطور مثال \$1 به (Nghbr) و \$3 به (Number) اشاره دارد که به ترتیب به نام محله و اجاره‌بها اشاره می‌کند. همانطور که دیده می‌شود مجموعه‌ی خروجی می‌تواند پایگاه داده‌ای از اطلاعات موجود در متن بوجود آورد.

روش تولید قواعد در WHISK مانند روشی است که در Wrapper Induction [۹] استفاده می‌شود. به‌اصفه‌ی اینکه از اطلاعات لغوی، گفتاری و نحوی موجود در متن برای قوی‌تر کردن قواعد استفاده می‌شود. سیستم WHISK، یادگیری قواعد را با کاربر انجام می‌دهد تا مثال‌هایی را برای مشخص کردن خروجی به کاربر ارائه کند که برای تصمیم‌گیری بعدی مفیدتر باشد و منجر به مجموعه قواعد استخراج کاراتر شود. برای تولید خروجی، قواعد بدست آمده بر اساس معیاری با هم مقایسه می‌شوند و باقاعدۀ مناسب‌تر انتخاب می‌شود. معیاری که در WHISK استفاده شده است میزان پوشش صحیحی است که قواعد بر روی مجموعه‌ی آموزشی دارند یعنی نسبت تعداد استخراج‌های درست بر کل استخراج‌هایی که بر روی مجموعه‌ی آموزشی داشته‌اند.

## ۳-۲ داده‌کاوی

هدف اصلی در داده‌کاوی کشف دانش شناخته نشده یا ضمنی در مجموعه‌های داده‌ای است.<sup>[۱۰]</sup> در داده‌کاوی یا کشف دانش از پایگاه داده (KDD)<sup>۳</sup> فرض می‌شود که داده در قالب پایگاه‌های داده رابطه‌ای می‌باشد. فرآیند عمومی در داده‌کاوی یا کشف دانش از پایگاه داده، در شکل ۲ آورده شده است.<sup>[۱۱]</sup>

تکنیک‌های مورد استفاده در داده‌کاوی در سه بخش تقسیم‌بندی می‌شوند<sup>[۱۰]</sup>: دسته‌بندی، خوشه‌بندی و ساختن قواعد تداعی‌گر.

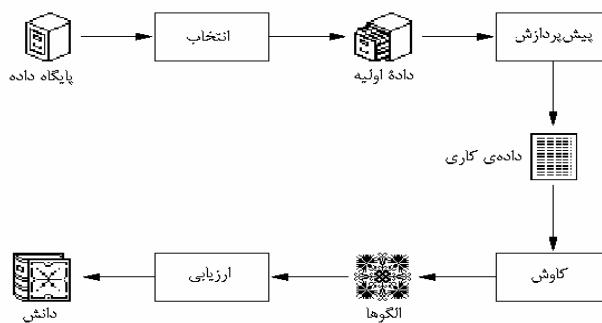
هدف در خوشه‌بندی پیشگویی دسته‌ی یک شیء جدید بر اساس مثال‌های آموزشی می‌باشد. بیشترین توجه‌ها در این قسمت بر روی الگوریتم‌های سازنده‌ی درخت‌های تصمیم و شبکه‌های عصبی می‌باشد. در خوشه‌بندی، مجموعه‌ای از اشیاء بر اساس شباهت مقادیر ویژگی‌هایشان خوشه‌بندی می‌شوند. یکی از کاربردهای مهم خوشه‌بندی ساخت رده‌بندی<sup>۴</sup> می‌باشد. در خوشه‌بندی اغلب از الگوریتم‌های خوشه‌بندی سلسله مراتبی جمع‌شونده<sup>۵</sup> استفاده می‌شود.

<sup>3</sup> Knowledge Discovery from Database

<sup>4</sup> Taxonomy

<sup>5</sup> Hierarchical Agglomerative Clustering

# فناوری اطلاعات و دانش



شکل ۲: فرآیند داده کاوی یا کشف دانش از پایگاه داده

در ساختن قواعد تداعی گر یا پیشگو فرم قواعد بصورت زیر می باشد:

$$A_1 \wedge A_2 \wedge \dots \wedge A_m \longrightarrow B_1 \wedge B_2 \wedge \dots \wedge B_n$$

هر یک از  $A_i$  ها و  $B_j$  ها مقادیر ویژگی ها یا ستون های یک جدول می باشد. این قاعده نشان می دهد که اگر رابطه  $A_1 \wedge A_2 \wedge \dots \wedge A_m$  میان مقادیر ستون های جدول برقرار باشد، می توان نتیجه گرفت که رابطه  $B_1 \wedge B_2 \wedge \dots \wedge B_n$  نیز برقرار است. الگوریتم های مختلفی برای ساخت این قواعد وجود دارد که از جمله می توان به RIPPER، A-Priori و C4.5RULES اشاره کرد [۱۲].

در این مقاله برای کارتر کردن سیستم استخراج اطلاعات، تکنیک ساختن قواعد تداعی گر را به آن اضافه خواهیم کرد که از این به بعد با نام قواعد پیشگو به آنها اشاره خواهیم کرد. برای این منظور از الگوریتم C4.5RULES استفاده می کنیم.<sup>6</sup>

در سال ۹۳، آقای Quinlan با گسترش الگوریتم استقرای ID3، روش C4.5 را پیشنهاد کرد. در ID3، پیشگویی دسته هر ورودی بر اساس مقادیر ویژگی هایش بود برای این منظور درخت تصمیمی بدست می آمد که هر یک از گره های میانی حاوی شرطی بر روی یک ویژگی بود و برگ های درخت، دسته ای را مشخص می کرد که ورودی به آن تعلق داشت. هر یک از مسیر های موجود از ریشه تا یکی از برگ ها نشان دهنده یک قاعده بود. پیمایش تمام مسیر های مذبور منجر به ایجاد یک مجموعه اولیه از قواعد پیشگو می شد. سپس در C4.5 مجموعه قواعد اولیه هرس می شود و شروط اضافی از بخش مقدم قواعد حذف و سپس قواعد تکراری حذف می شوند. خروجی C4.5، مجموعه قواعد به مراتب کوچکتر و ساده تر نسبت به مجموعه قواعد اولیه می شد. C4.5RULES نسخه بهبود یافته نرم افزاری C4.5 می باشد که یکی از پر کاربرد ترین برنامه ها در داده کاوی به حساب می آید. پس از اعمال الگوریتم C4.5RULES بر روی یک پایگاه داده قواعدی بصورت زیر خواهیم داشت [۴].

$$A_1 \wedge A_2 \wedge \dots \wedge A_m \longrightarrow B$$

هر یک از  $A_i$  ها و  $B$  می تواند بصورت  $a \in Col_i$  یا  $a \notin Col_i$  باشد که نشان دهنده این است که مقدار ویژگی ستون  $i$  آن برابر  $a$  باشد یا خیر. بر حسب اینکه  $B$  کدام یک از این حالات را دارد، بترتیب قواعد پیشگو مثبت و منفی خواهیم داشت.

### ۳- استفاده از C4.5RULES برای بهبود کارآیی سیستم های استخراج اطلاعات

برای بکارگیری تکنیک های داده کاوی، خروجی سیستم استخراج اطلاعات باید به صورت پایگاه داده رابطه ای باشد. خروجی سیستم استخراج اطلاعات بصورت یک ساختار با چند قاب می باشد که به راحتی تبدیل به جدول با چند ستون خواهد شد.

سیستم استخراج اطلاعات در یک دامنه خاص فعالیت می کند. دامنه ای که به عنوان مثال در این مقاله در نظر گرفته شده است، دامنه ای آگهی های استخدامی مربوط به رشته کامپیوتر می باشد. این آگهی ها را می توان در گروه های خبری

<sup>6</sup> <http://www.cs.uregina.ca/~dbd/cs831/notes/ml/dtrees/c4.5/c4.5rules.html>

# فناوری اطلاعات و دانش

تهران / دانشگاه صنعتی امیرکبیر / ۳-۵ خرداد ۱۳۸۴

اینترنتی یافت. سیستم استخراج اطلاعات از روی این آگهی‌ها، خروجی‌هایی مانند نمونه‌ای که در شکل ۳ آورده شده است، استخراج می‌کند.

نکته‌ی دیگر این است که بعضی از مقادیر گرچه از نظر لغوی با هم متفاوت هستند ولی از نظر ماهیت یکسانند. بطور مثال مقادیر VB و Visual Basic Microsoft هر سه به زبان برنامه‌نویسی ویژوال بیسیک اشاره دارند. در این موقع لازم است که به عنوان یک پیش‌پردازش، یکسان‌سازی مقادیر صورت گیرد. مرحله‌ی بعدی باین‌ری کردن محتویات جدول است بطوری که مقادیر فقط صفر یا یک باشند<sup>[۴]</sup>. برای این منظور به جای هر یک از ستون‌های انتخابی، چند ستون به تعداد مقادیر ممکن آن قرار می‌دهیم. حال به ازای هر یک از مقادیر موجود در ستون پیشین، محتویات ستون همنام با آن مقادیر را یک قرار می‌دهیم و بقیه‌ی ستون‌ها را با صفر پر می‌کنیم. اینک جدول آماده اعمال الگوریتم C4.5RULES است. این الگوریتم در ابتدا با ساخت درخت‌های تصمیم و سپس هرس کردن و ترجمه‌ی آن‌ها مجموعه‌ای از قواعد بدست می‌آورد که وجود روابطی بین محتویات ستون‌های جدول را پیشگویی می‌کنند. هر چه مجموعه‌ی آموزشی بزرگتری داشته باشیم، مجموعه قواعد کاراتری خواهیم داشت. در شکل ۴ نمونه‌هایی از قواعد پیشگوی بدست آمده نشان داده شده است.

Title: DBMS Administrator
City: Austin
Area: EC, DB Applications
Lanquage: VB, VC, C++, Delphi, Java
Platform: Windows, Unix
Application: SQL Server, Oracle
Desires years of experience: 5
Required degree: BS

شکل ۳: نمونه خروجی سیستم استخراج اطلاعات (دامنه‌ی آگهی‌های استخدامی رشته‌ی کامپیوتر)

SQL Server  $\in$  Application  $\wedge$  DB  $\in$  Application Area  $\longrightarrow$  Oracle  $\in$  Application  
 VC++  $\in$  Language  $\wedge$  VB  $\in$  Language  $\longrightarrow$  Windows  $\in$  Platform  
 $\neg$ (Windows  $\in$  Platform)  $\wedge$  Unix  $\in$  Platform  $\longrightarrow$   $\neg$ (VC  $\in$  Language)

شکل ۴: نمونه قواعد بدست آمده از C4.5RULES

در شکل بالا قاعده‌ی اول و دوم قواعد پیشگوی مثبت و قاعده‌ی سوم قاعده‌ی پیشگوی منفی می‌باشد. بطور مثال قاعده آخر، پیشگویی می‌کند که اگر در مهارت‌های خواسته شده در آگهی، سیستم عامل ویندوز وجود نداشت و در ضمن سیستم عامل یونیکس وجود داشت می‌توان نتیجه‌گیری کرد که احتمالاً زبان ویژوال سی جزو مهارت‌های درخواستی نخواهد بود. در مرحله‌ی بعد می‌خواهیم از این نوع قواعد برای بهبود کارآیی سیستم‌های استخراج اطلاعات استفاده کنیم. کارآیی سیستم‌های استخراج اطلاعات با دو معیار Precision و Recall اندازه‌گیری می‌شوند که به صورت زیر تعریف می‌شوند<sup>[۲]</sup>:

$$\text{Recall} = \frac{\text{تعداد کل استخراج‌جهای ممکن در متن}}{\text{تعداد استخراج‌جهای درست}}$$

$$\text{Precision} = \frac{\text{تعداد کل استخراج‌جهای انجام شده}}{\text{تعداد استخراج‌جهای درست}}$$

همچنین برای ترکیب این دو، معیار F-Measure با تعریف زیر استفاده می‌شود:

$$\text{F-Measure} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

معمولًا در سیستم‌های استخراج اطلاعات مقدار Recall بسیار کمتر از مقدار Precision است<sup>[2,13]</sup>. با توجه به تعاریف آن‌ها مشخص است که مقادیر آن‌ها با هم مرتبط‌اند. اغلب روش‌هایی که برای بالا بردن یک معیار استفاده می‌شود باعث کاهش معیار دوم می‌شود و در نتیجه معیار F-Measure تغییر چندانی نمی‌کند. اما ما با استفاده از قواعد پیشگوی مقادیر می‌توانیم Precision و Recall را بدون ضربه زدن به یکدیگر بهبود داده‌ایم. در زیر چند پیشنهاد برای استفاده از قواعد پیشگوی آمده است و نتایج اعمال آنها بر روی سیستم WHISK در بخش ارزیابی بیان شده است.

# فناوری اطلاعات و دانش

تهران / دانشگاه صنعتی امیرکبیر / ۳-۵ خرداد ۱۳۸۴

۱- خود قواعد بدست آمده به عنوان یک خروجی بالرژش در کنار خروجی‌های بدست آمده از سیستم استخراج اطلاعات قابل استفاده است. بطور مثال از روی قواعد بدست آمده در دامنه‌ی مذبور می‌توان به یک فرد توصیه کرد که اگر می‌خواهد در زمینه‌ی گرافیک در حوزه‌ی نوشن بازی‌های کامپیوترا دارای بازار کار خوبی داشته باشد چه مهارت‌هایی را بیشتر مد نظر قرار دهد.

۲- قواعد پیشگویی مثبت: با استفاده از قواعد پیشگویی مثبت می‌توان پیشگویی کرد که اگر شرطی خاص بین مقادیر خروجی برقرار باشد باید موردی خاص در خروجی وجود داشته باشد. اگر اطلاعات پیشگویی شده، استخراج نشده باشد، می‌توان آن‌ها را به خروجی اضافه کرد. بدین وسیله معیار Recall بدون صدمه زدن به Precision افزایش یافته است. بطور مثال فرض کنید که قاعده  $VC \in language \rightarrow Win \in platform$  برخوان سیستم عامل استخراج نشده است. پیشگویی می‌شود که احتمالاً Win در برنامه‌نویسی استخراج شده است ولی Win در قاب  $language$  با صفر پر شده است و هدف ما افزایش معیار Recall و استخراج بیشتر Win در قاب platform جاافتاده است. برای اینکه گواهی برای این مدعای داشته باشیم در متن ورودی جستجو کرده و اگر کلمه‌ی Win پیدا شد، فرض می‌شود که پیشگویی صحیح بوده است و آن مورد را به خروجی اضافه می‌کنیم.

ممکن است که چند قاعده پیشگوی وجود داشته باشد و بکارگیری آن‌ها با ترتیب‌های مختلف منجر به نتایج متفاوت شود. سؤال این است که چه ترتیبی بهتر است. پیشنهاد می‌کنیم قواعدی که دارای عبارت‌های نقیض دار بیشتری می‌باشند در الیت قرار بگیرند. زیرا می‌دانیم که بیشتر محتویات جداول با صفر پر شده است و هدف ما افزایش معیار Recall و استخراج بیشتر اطلاعات است.

۳- قواعد پیشگوی منفی: قواعد بدست آمده از C4.5RULES دارای درصد احتمال درستی هستند که به صورت آماری از روی پایگاه داده بدست آمده‌اند. در مورد قبل که از قواعد پیشگویی مثبت استفاده کردیم، می‌توانستیم شاهدی از متن برای تصدیق پیشگویی قاعده بیاوریم. اما درباره قواعد پیشگویی منفی چنین تصدیقی وجود ندارد. بطور مثال قاعده‌ی  $(VC \in area \rightarrow (Win \in area \rightarrow VC \in language))$  پیشگویی می‌کند که در صورتیکه در خروجی Win در قاب area نباشد، نباید VC نیز در قسمت language آمده باشد. اگر با درست بودن فرض قاعده، در خروجی VC در قاب language آمده باشد بدین معنی است که در متن وجود داشته است و نمی‌توان تصدیقی برای پیشگویی قاعده‌ی مذبور از درون متن بدست آورد. دو پیشنهاد برای استفاده از قواعد پیشگویی منفی داریم: اول اینکه، از قواعد با درصد احتمال درستی بالا استفاده کنیم و پیشگویی آن‌ها را با تکیه به بالا بودن احتمال درستی‌شان بذیریم و اگر در خروجی موردنی برخلاف پیشگویی داشته باشیم، آن را تصحیح کنیم. در این روش تعداد پیشگویی‌ها خیلی کم می‌شود ولی در عین حال، پیشگویی، اکثر اوقات صحیح خواهد بود.

در روش دوم از قواعد پیشگویی منفی با درصد احتمال درستی مناسب استفاده می‌کنیم. وقتی متنی به سیستم استخراج اطلاعات داده می‌شود، اکثر اوقات بیش از یک قاعده‌ی استخراج اطلاعات می‌تواند بر روی آن بکار برد شوند. به عبارت دیگر چون هر قاعده یک الگو را در متن جستجو می‌کند، ممکن است متن ورودی با چند الگو مطابقت کند. برای انتخاب مناسب‌ترین قاعده، علاوه بر معیار مقایسه‌ی قواعد می‌توان از قواعد پیشگویی منفی استفاده کرد و امتیاز بیشتری برای قاعده‌ای در نظر گرفت که با پیشگویی‌های انجام شده مطابقت داشته باشد. با اعمال این روش صحت درستی استخراج‌ها بیشتر خواهد شد. به عبارت دیگر Precision را بهبود داده‌ایم.

البته می‌توان از تلفیق این دو روش نیز استفاده کرد. بدین طریق که حدی را تعیین کرد و از قواعد با احتمال درستی بالاتر از آن مستقیم در تصحیح خروجی استفاده کرد و از بقیه قواعد برای انتخاب قاعده‌ی استخراج بهینه استفاده کرد. در قسمت ارزیابی از روش دوم استفاده شده است.

۴- قواعد پیشگویی نحوی: سیستم‌های استخراج اطلاعات که برای متن دون ساختار طراحی شده‌اند بر روی متن ورودی پردازش نحوی انجام می‌دهند و از اطلاعات نحوی اضافه شده به متن برای قدرتمند کردن قواعد استخراج اطلاعات استفاده می‌کنند. در خروجی چنین سیستم‌هایی نقش نحوی مقادیر ظاهر شده در خروجی مشخص است و می‌توان پایگاه داده‌ای از نقش‌های مقادیر خروجی ساخت. چون تعداد نقش‌های نحوی محدود و شمارش‌پذیر هستند این پایگاه داده را می‌توان برای اعمال الگوریتم C4.5RULES استفاده کرد و روابط جدید نحوی بین مقادیر سازنده‌ی خروجی‌ها کشف کرد. روش کار در اینجا

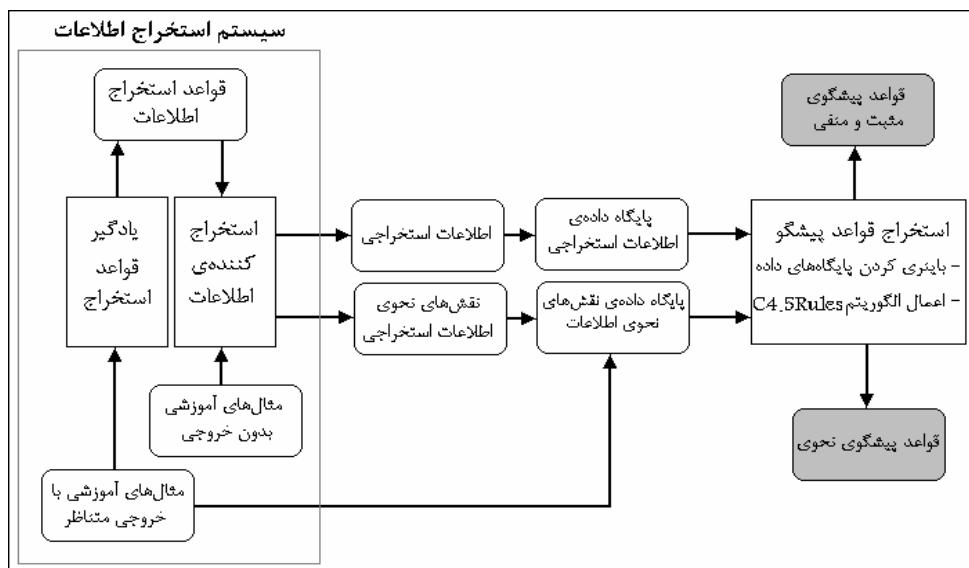
# فناوری اطلاعات و دانش

تهران / دانشگاه صنعتی امیرکبیر / ۳-۵ خرداد ۱۳۸۴

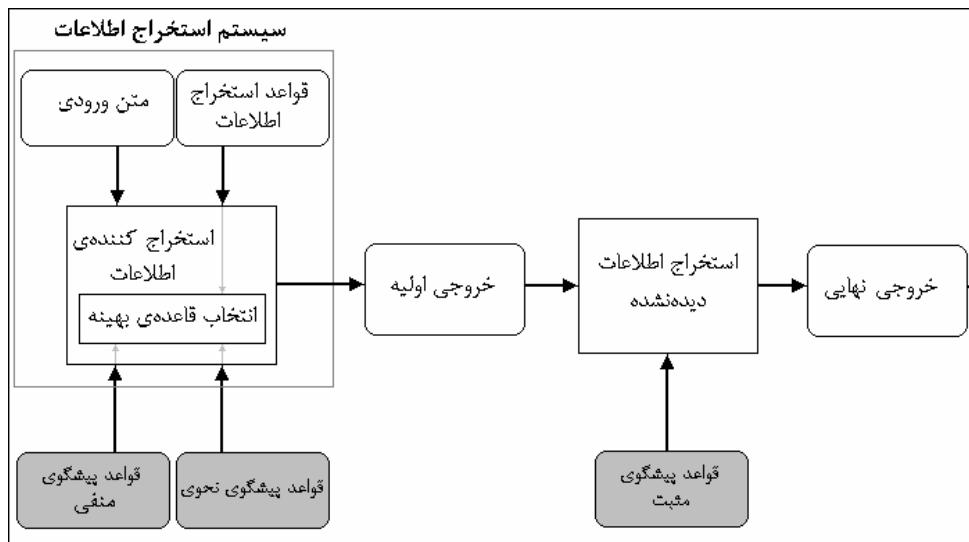
مانند مورد قبل است با این تفاوت که به جای پایگاه داده‌ی خروجی از پایگاه داده‌ی نقش‌های نحوی خروجی‌ها استفاده می‌کنیم. حاصل این کار بدست آوردن قواعدی است که وجود روابطی بین نقش‌های نحوی قاب‌های خروجی را پیشگویی می‌کنند.

نحوی استفاده از این قواعد نیز مانند مورد قبل است و برای انتخاب قاعده‌ی استخراج بهینه در هنگام تولید خروجی استفاده می‌کنیم. با این کار بهبودی دیگر در معیار Precision حاصل می‌شود. در شکل ۵ معماری پیشنهادی سیستم یادگیرنده‌ی قواعد پیشگو نشان داده شده است. این کار با استفاده از خروجی‌های بدست آمده از سیستم استخراج اطلاعات و مثال‌های آموزشی (که قبل از خروجی‌ها ایشان به صورت دستی بدست آمده است) صورت می‌گیرد.

در شکل ۶ معماری سیستم استخراج اطلاعات در حالت استفاده از قواعد پیشگو طراحی شده است. از قواعد پیشگوی منفی و قواعد پیشگوی نحوی برای انتخاب قاعده‌ی بهینه برای تولید خروجی اولیه استفاده می‌شود و سپس از قواعد پیشگوی مثبت برای اضافه کردن اطلاعات جافتاده‌ی احتمالی استفاده می‌شود.



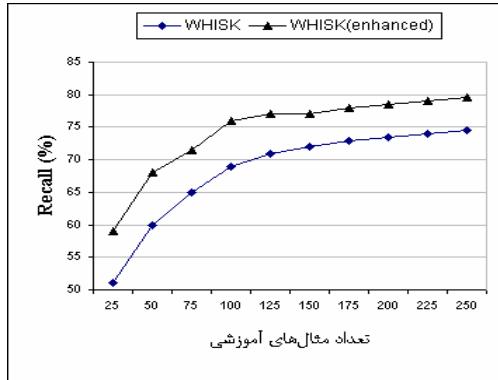
شکل ۵: معماری سیستم یادگیر قواعد پیشگو



شکل ۶: نحوی استفاده از قواعد پیشگو در سیستم استخراج اطلاعات

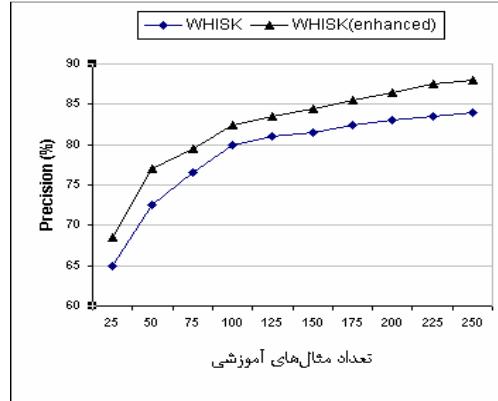
## ۴- ارزیابی

سیستم استخراج اطلاعاتی که برای ارزیابی استفاده شد، سیستم WHISK می‌باشد. کد برنامه و توضیحات آن در اینترنت موجود است.<sup>7</sup> دامنهٔ مورد نظر، آگهی‌های استخدام شغل‌های کامپیوتری می‌باشد. مثال‌های مورد نیاز از گروه‌های خبری اینترنتی مانند www.austinjobs.com جمع‌آوری شد. در مجموع ۳۰۰ مثال بدست آمد. سپس خروجی آن‌ها به صورت دستی محاسبه شد. از ۲۵۰ عدد آن‌ها برای آموزش سیستم WHISK و یادگیری قواعد استخراج استفاده شد. از همان مثال‌ها برای بدست آوردن قواعد پیشگو استفاده شد. به عبارت دیگر از مثال‌های آموزشی بدون خروجی که در شکل ۴ آمده است برای یادگیری قواعد استفاده نشده است. چون در عمل دیده شد زمانیکه از این مثال‌ها استفاده می‌شد، معیارهای کارآیی پایین‌تر می‌آمد.



شکل ۷: مقایسهٔ معیار Recall

سپس از ۵۰ عدد مثال‌های قبلی به علاوهٔ ۵۰ مثال باقیمانده برای تست سیستم و بدست آوردن معیارهای کارآیی استفاده شده است. برای یادگیری قواعد پیشگو از فیلد‌های area, application, platform, language و استفاده شده است. چون ارتباط بیشتری با هم داشتند.



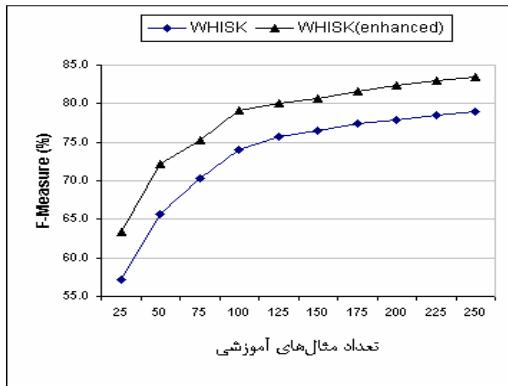
شکل ۸: مقایسهٔ معیار Precision

اشکال ۷، ۸ و ۹ مقایسه مقادیر معیارهای کارآیی را نشان می‌دهد. مقایسه، توسط مقادیر معیارهای سه‌گانه‌ی ذکر شده در قبل صورت گرفته است. منحنی پایینی برای زمانی است که از سیستم WHISK به تنهاًی استفاده شده است و منحنی بالایی برای زمانی است که تکنیک استفاده از قواعد پیشگو را به آن اضافه کرده‌ایم. هر چقدر مثال‌های آموزشی بیشتر شده است، معیارها بهبود بیشتری را نشان داده‌اند ولی از یک مرحله به بعد اضافه کردن آن‌ها بهبود چندانی حاصل نکرده است. بطور متوسط معیارهای Precision, Recall و F\_Measure نسبت به قبیل (استفاده از WHISK به تنهاًی) بترتیب ۶٪، ۴٪ و

<sup>7</sup> <http://www.cs.washington.edu/homes/soderlan/WHISK.tar.gz>

# فناوری اطلاعات و دانش

۵٪ افزایش یافته است. با توجه به اینکه مقادیر معیارهای کارآیی سیستم WHISK در مقایسه با سیستم‌های دیگر استخراج اطلاعات در حد قابل قبولی قرار دارد [۱]، نتایج بدست آمده در جهت بهبود کارآیی سیستم‌های استخراج اطلاعات امیدوار کننده می‌باشد.



شکل ۹: مقایسه معيار F-Measure

## ۵- کارهای انجام شده

با توجه اینکه موضوع استخراج اطلاعات از متن در سال‌های اخیر مورد توجه بیشتر قرار گرفته است در گذشته کار چندانی بر روی استفاده از تکنیک‌های داده‌کاوی در استخراج اطلاعات از متن صورت نگرفته است. در سال ۹۵ آقایان Feldnom و Dagan سیستمی را با نام KDT ارائه کردند [۱۴] که از پایگاه داده‌ی متنی، دانش استخراج می‌کرد. سیستم مزبور متنون را به صورت دستی و غیر خودکار دسته‌بندی می‌کرد. به عبارت دیگر KDT از یک دسته‌بندی کننده‌ی خودکار یا سیستم استخراج اطلاعات استفاده نمی‌کرد و در مقاله‌ی مربوطه نیز هیچ بحثی درباره‌ی استفاده از دانش استخراجی برای بهبود عملکرد سیستم نشده است. در سال‌های اخیر توجهات زیادی به مبحث متن‌کاوی جلب شده است. با این حال سیستم‌های عملیاتی و ارزیابی‌های تجربی کمی وجود دارد.

## ۶- زمینه‌های تحقیقاتی و نتیجه‌گیری

یکی از مراحلی که برای برای تبدیل خروجی سیستم‌های استخراج اطلاعات به پایگاه‌های داده‌ی مورد قبول الگوریتم C4.5RULES وجود داشت، یکپارچه سازی عبارات هم‌معنی بود. این کار در پیاده‌سازی صورت گرفته به صورت دستی انجام شد ولی می‌توان بر روی استفاده از روش‌های خودکار مانند محاسبه‌ی شباهت دو عبارت از نظر حروف تشکیل دهنده [۱۵] تحقیق کرد.

در بررسی صورت گرفته مقادیر قاب‌های خروجی همیشه گسته فرض شد و برای حالتی که مقادیر پیوسته باشند مانند طول یا قیمت بحثی نشد. یافتن روش‌هایی برای پیشگوی مقادیر پیوسته می‌تواند موضوعی برای تحقیقات آتی باشد. در این مقاله، اینکه از چه فیلدهایی از جدول برای تولید قواعد پیشگو استفاده شود به صورت ذهنی بدست آمد. اگر بتوان ارتباط میان ستون‌های جدول را از روی محتويات آن‌ها پیدا کرد، می‌توان انتخاب فیلدهای مناسب برای یادگیری قواعد پیشگو را خودکار انجام داد.

مطمئناً در آینده تحقیقات بیشتری بر روی استفاده از تکنیک‌های موجود در موضوعات دیگر هوش مصنوعی برای بهبود کارآیی سیستم‌های استخراج اطلاعات صورت خواهد گرفت. نتیجه‌گیری اینکه، همکاری استخراج اطلاعات از متن و داده‌کاوی می‌تواند برای کاربردهای مختلفی مفید باشد. سیستم استخراج اطلاعات خروجی‌هایی را در اختیار قرار می‌دهد که داده‌کاوی می‌تواند از روی آن، دانش مفیدی برای کاربردهای مختلف بویژه متن‌کاوی تولید کند.

# فناوری اطلاعات و دانش

تهران / دانشگاه صنعتی امیرکبیر / ۳-۵ خرداد ۱۳۸۴

متن کاوی موضوع نسبتاً جدیدی است که در آن پردازش‌های زبان طبیعی، یادگیری ماشین و بازیابی اطلاعات با هم همکاری دارند [۱۶]. ادغام این روش‌ها می‌تواند منجر به ایجاد تکنیک‌های جدید و مفید برای استخراج دانش از مجموعه‌های حجیم متنی شوند.

## مراجع

- [1] S. Soderland, *Learning Information Extraction Rules for Semi-Structured and Free Text*, Kluwer Academic Publishers, Boston, 1999.
- [2] DARPA., ed. Proceedings of the Fifth DARPA Message Understanding Conference, San Mateo, CA: Morgan Kaufman, 1995.
- [3] C. Cardie, and R. J. Mooney, *Machine learning and natural language (introduction to special issue on natural language learning)*, Machine Learning Magazine, No 34, pp 5-9, 1999.
- [4] J. R. Quinlan, *C4.5: Programs for Machine Learning*, San Mateo, CA: Morgan Kaufman.
- [5] F. Peng, *Models Development in IE Tasks*, CS685 Project, Computer Science Department, University of Waterloo, <http://ai2.uwaterloo.ca/~f3peng>, 2001.
- [6] M. E. Callif, and R. J. Mooney, *Relational learning of pattern-match rules for information extraction*, In Proceedings of the Sixteenth National Conference on Artificial Intelligence, pp 328-324, 1999.
- [7] H. Ahonen-myka, *Information Extraction from text*, Information Extraction Course Slides, Carnegie Mellon School of Computer Science, spring 2002.
- [۸] ا. عبداله زاده، و. شکری. روش‌های "مهندسی دانش" و "سیستمهای یادگیر" برای طراحی سیستمهای استخراج اطلاعات از متن، مجموعه مقالات چهارمین کنفرانس دانشجویی انجمن کامپیوتر ایران، نجف آباد. ۱۳۸۱.
- [9] N. Kushmerick, D. Weld, and R. Doorenbos, *Wrapper induction for information extraction*, Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence, pp 729-737, 1997.
- [10] M. Holsheimer, M. Kerstern, H. Mannila, and H. Toivonen, *A perspective on databases and data mining*, Proceedings of the First International Conference on Knowledge Discovery and Data Mining, AAAI Press, PP 150-155, 1995.
- [11] G. J. Williams, and Z. Huang, *Modeling the KDD Process, A Four Stage Process and Four Element Model*, CSIRO Division of Information Technology, 1996.
- [12] R. Agrawal, and R. Sirkant, *Fast algorithms for mining association rules in large databases*, VLDB'94, 1994.
- [13] DARPA., ed. Proceedings of the Fifth DARPA Message Understanding Evaluation and Conference, San Mateo, CA: Morgan Kaufman, 1993.
- [14] R. Feldman, and I. Dagan, *Knowledge discovery in textual database (KDT)*, In Proceedings of the First International Conference on Knowledge Discovery and Data Mining, 1995.
- [15] E. S. Ristad, and P. N. Yianilos, *Learning string edit distance*, IEEE Transactions on Pattern Analysis and Machine Intelligence 20(5), 1998.
- [16] M. Hearst, *Untangling text data mining*, In Proceedings of the 37<sup>th</sup> Annual Meeting of the Association for Computational Linguistic, pp 3-10, 1999.