



1. 1180

۸۷/۱/۱۰۵۸۷۵  
-----  
۸۷/۱۲/۴



دانشکده مهندسی

پایان نامه کارشناسی ارشد در رشته ی مهندسی کامپیوتر- نرم افزار

استخراج مجموعه آیتم های فراوانی در محیط توزیع شده

بوسیله ی:

الهام پرنیان

استاد راهنما:

دکتر محمد هادی صدرالدینی

کتابخانه مرکزی  
شیراز

۱۳۸۷ ۸۷ ۸

شهریورماه ۱۳۸۷

۱۰۸۸۵۵

به نام خدا

## استخراج مجموعه آیتم های فراوانی در محیط توزیع شده

به وسیله  
الهام پرنیان

### پایان نامه

ارائه شده به تحصیلات تکمیلی دانشگاه به عنوان بخشی  
از فعالیت های تحصیلی لازم برای اخذ درجه کارشناسی ارشد

در رشته‌ی

مهندسی کامپیوتر (نرم افزار)

از دانشگاه شیراز

شیراز

جمهوری اسلامی ایران

ارزیابی شده توسط کمیته پایان نامه با درجه: عالی

دکتر محمد هادی صدرالدینی، استادیار بخش مهندسی و علوم کامپیوتر (رییس کمیته).....  
دکتر غلامحسین دستغیبی فرد، استادیار بخش مهندسی و علوم کامپیوتر.....  
دکتر فریبرز سبحان منش، استادیار بخش مهندسی و علوم کامپیوتر.....

شهریور ۱۳۸۷

تقدیم به تمامی آموزگاران الفبای زندگی

به مادر، باران فداکاری

و به پدر، آسمان شکیبایی

## سپاسگزاری

اکنون که به یاری خداوند متعال موفق به اتمام این پایان نامه شده‌ام، از زحمات اساتید بزرگوار جناب آقای دکتر محمد هادی صدرالدینی، دکتر غلامحسین دستغیبی، فرد و دکتر فریبرز سبحان منش سپاسگزارم.

از تمامی دوستان و خانواده عزیزم که مشوقم بودند و همچنین از تمامی کادر بخش مهندسی و علوم کامپیوتر که در انجام این کار با اینجانب همکاری داشته‌اند، صمیمانه سپاسگزارم. و از مرکز تحقیقات مخابرات ایران به دلیل حمایت مالی از این پایان نامه، کمال تشکر را دارم.

## چکیده

استخراج مجموعه آیت‌های فراوانی در محیط توزیع شده

بوسیله‌ی:

الهام پرنیان

حجم روزافزون داده‌ها در فایل‌ها، پایگاه‌های داده و دیگر مخازن داده‌ای، توسعه ابزارهای تجزیه تحلیل و تفسیر داده‌ها و استخراج اطلاعات جالب از این داده‌ها را ایجاب می‌نماید. این اطلاعات در پروسه‌های تصمیم‌گیری سازمانها مورد استفاده قرار می‌گیرند. داده‌کاوی، کشف الگوهای پنهان و اطلاعات مفید از پایگاه داده‌هاست و یکی از قدم‌های مهم در پروسه کشف دانش است. ارتباطات پنهان بین داده‌ها در بانک اطلاعاتی تراکنشی، قوانین وابستگی نامیده می‌شود. کشف قوانین وابستگی یکی از مهمترین شاخه‌های داده‌کاوی است که زمینه‌های کاربردی فراوانی دارد. در ابتدا الگوریتم‌های زیادی در زمینه کشف این قواعد به صورت مرکزی، ارائه شد. اما با ظهور سازمانهایی که به صورت جغرافیایی توزیع شده هستند و با توجه به حجم بالای ارتباطات شبکه‌ای و لزوم حفظ محرمانگی داده‌ها، کاوش قوانین وابستگی در محیط‌های توزیع شده هم مد نظر قرار گرفت. کاوش مجموعه اقلام فراوان اولین قدم در کشف قوانین وابستگی است. به طور معمول این کار باعث ایجاد یک مجموعه بزرگ از مجموعه قلم‌ها می‌شود که تعداد زیادی از آنها اضافی هستند. در سال‌های اخیر، روشهایی در زمینه کاوش مجموعه قلم‌های تکرار شونده بسته پیشنهاد شده است. این مجموعه قلم‌ها، یک نمایش خلاصه و کامل از مجموعه قلم‌های تکرار شونده در پایگاه داده ارائه می‌دهند. حجم ارتباطات و زمان لازم برای اجرا، از پارامترهای مهم در کشف قوانین وابستگی در محیط توزیع شده هستند. از آنجایی که ارسال مجموعه اقلام فراوان در یک محیط توزیع شده، حجم بالایی از ارتباطات شبکه‌ای را سبب می‌شود، ما در این پایان‌نامه با استفاده از تکنیک فشرده‌سازی مجموعه اقلام در یکی از معروفترین الگوریتم‌های توزیع شده به نام او. دی. ای. ام، یک روش جدید ایجاد کرده ایم و آن را با داده‌های واقعی مورد آزمایش قرار داده‌ایم. با استفاده از روش جدید حجم ارتباطات در الگوریتم او. دی. ای. ام به نحو قابل قبولی کاهش می‌یابد. نتایج ارزیابی صورت گرفته بر روی دو الگوریتم او. دی. ای. ام و الگوریتم جدید، کارایی برتر الگوریتم جدید را از دید حجم ارتباطات شبکه‌ای و زمان اجرا نشان می‌دهد. این روش را می‌توان در هر الگوریتم توزیع شده به کار برد.

## فهرست مطالب

صفحه	عنوان
	فصل اول: داده کاوی و کاوش قواعد وابستگی
۲	۱-۱ مقدمه
۳	۲-۱ داده کاوی
۶	۳-۱ مراحل داده کاوی
۷	۴-۱ روش های داده کاوی
۷	۵-۱ کاوش قواعد وابستگی
۷	۶-۱ تعریف مسئله کاوش قواعد وابستگی
۹	۷-۱ الگوریتم ای. پریوری
۱۲	۱-۷-۱ تابع Subset
۱۴	۸-۱ بهینه سازی های انجام شده بر روی الگوریتم ای. پریوری
۱۴	۱-۸-۱ الگوریتم های ای. پریوری هیبرید و ای. پریوری تی. آی. دی
۱۵	۹-۱ نمونه گیری
۱۶	۱۰-۱ پارتیشن بندی
۱۶	۱۱-۱ روش اف. پی. گرات
۱۶	۱۱-۱ هش کردن و هرس کردن مستقیم
۱۶	۱۲-۱ مجموعه اقلام غیرقابل اشتقاق
۱۶	۱-۱۲-۱ قوانین استنتاجی
۱۹	۲-۱۲-۱ الگوریتم ان. دی. آی
۲۱	۳-۱۲-۱ تولید همهی مجموعه اقلام با استفاده از مجموعه اقلام غیرقابل اشتقاق
۲۴	۱۳-۱ ساختمان داده Trie

۲۴	۱-۱۳-۱ روشهای شمارش پشتیبانی مجموعه اقلام در Trie
۲۵	۱-۱۳-۱ ذخیره سازی داده ورودی

فصل دوم : کاوش مجموعه قلم های تکرار شونده بسته

۲۷	۱-۲ کاوش مجموعه قلم های تکرار شونده بسته
۲۷	۲-۲ ضرورت کاوش مجموعه قلم های تکرار شونده بسته
۲۹	۳-۲ الگوریتم های کاوش مجموعه قلم های تکرار شونده بسته
۲۹	۱-۳-۲ الگوریتم ای کلوز
۳۲	۲-۳-۲ الگوریتم کلوزت
۳۳	۳-۳-۲ الگوریتم کلوزت پلاس
۳۳	۴-۳-۲ الگوریتم چارم
۳۴	۱-۴-۳-۲ شرح الگوریتم چارم
۳۶	۵-۳-۲ الگوریتم اف پی کلوز
۳۷	۱-۵-۳-۲ درخت سی اف آی
۳۸	۲-۵-۳-۲ تکنیک اف پی اری
۳۹	۳-۵-۳-۲ پیاده سازی الگوریتم اف پی کلوز
۴۱	۶-۳-۲ الگوریتم دی سی آی کلوز
۴۲	۱-۶-۳-۲ مجموعه اقلام بسته
۴۵	۲-۶-۳-۲ حذف بهینه تکراری ها
۴۶	۳-۶-۳-۲ طرح الگوریتم دی سی آی کلوز
۴۷	۷-۳-۲ مجموعه اقلام بسته غیر قابل اشتقاق

فصل سوم : داده کاوی توزیع شده و کاوش توزیع شده قواعد وابستگی

۵۲	۱-۳ داده کاوی توزیع شده
۵۳	۲-۳ ضرورت داده کاوی توزیع شده
۵۴	۳-۳ موارد مهم در داده کاوی توزیع شده
۵۷	۴-۳ الگوریتم های توزیع شده کاوش قواعد وابستگی
۵۷	۱-۴-۳ الگوریتم توزیع شمارش
۵۹	۲-۴-۳ الگوریتم اف .دی .ام



۶۱	۳-۴-۳ الگوریتم او.دی.ای.ام.
۶۴	۴-۴-۳ الگوریتم دی.دی.ام.
۶۵	۵-۴-۳ کاوش مجموعه آیتمهای فراوان بسته توزیع شده
۶۵	۳-۴-۵-۱ کاوش مجموعه آیتم های بسته
۶۷	۳-۴-۵-۲ محاسبه فراوانی مجموعه آیتمهای بسته
۶۸	۳-۴-۶ کاوش مجموعه اقلام فراوانی در سکوهای نامتجانس

#### فصل چهارم: الگوریتمی جدید برای کاوش قواعد وابستگی در محیط توزیع شده

۷۰	۴-۱ نمایش بسته مجموعه آیتم ها
۷۰	۴-۲ مجموعه اقلام غیر قابل اشتقاق
۷۱	۴-۳ مجموعه اقلام بسته غیر قابل اشتقاق
۷۱	۴-۴ مجموعه اقلام بسته توزیع شده
۷۲	۴-۵ الگوریتم او دی ای ام
۷۲	۴-۶ الگوریتم سی ان دی آی، او دی ای ام
۷۳	۴-۶-۱ تابع delete_infrequent
۷۴	۴-۶-۲ تابع Superset_Support
۷۶	۴-۶-۳ پیاده سازی ای پریوری با استفاده از درخت Trie
۸۰	۴-۶-۴ به دست آوردن تمام اقلام فراوان از مجموعه اقلام بسته غیرقابل اشتقاق

#### فصل پنجم: ارزیابی کارایی و نتیجه گیری

۸۳	۵-۱ محیط آزمایش
۸۳	۵-۲ پایگاه داده های استفاده شده برای ارزیابی کارایی
۸۵	۵-۳ زمان اجرا
۸۹	۵-۴ تبادل پیام
۹۰	۵-۵ نتیجه گیری
۹۱	۵-۶ زمینه های تحقیقاتی آینده داده کاوی توزیع شده

۹۲	فهرست منابع
----	-------------

## فهرست شکل ها

صفحه	عنوان
۱۰	شکل ۱-۱ الگوریتم ای. پریوری
۱۱	شکل ۲-۱ مجموعه اقلام کاندید در الگوریتم ای. پریوری
۱۳	شکل ۳-۱ شمارش پشتیبانی مجموعه اقلام کاندید به کمک درخت هش
۱۷	شکل ۴-۱ کلیه قوانین استنتاجی ممکن برای مجموعه قلم abcd
۱۹	شکل ۵-۱ الگوریتم ان. دی. آی
۳۲	شکل ۱-۲ الگوریتم ای کلوز
۳۵	شکل ۲-۲ الگوریتم چارم- ترتیب لغت نامه‌ای
۳۷	شکل ۳-۲ اف پی تری
۳۸	شکل ۴-۲ ایجاد درخت مجموعه قلم های تکرار شونده بسته
۴۰	شکل ۵-۲ الگوریتم FP_Close
۴۲	شکل ۶-۲ شبکه تمام مجموعه اقلام فراوان
۴۹	شکل ۷-۲ مجموعه اقلام بسته، غیرقابل اشتقاق و بسته غیرقابل اشتقاق و فراوانی آنها
۵۰	شکل ۸-۲ شبه کد الگوریتم CNDI
۵۳	شکل ۱-۳ معماری کلی داده کاوی توزیع شده
۵۸	شکل ۲-۳ تکرار دوم از الگوریتم توزیع شمارش در یک سیستم توزیع شده دارای ۳ سایت
۷۴	شکل ۱-۴ تابع Delete_infrequent
۷۵	شکل ۲-۴ تابع Superset_support
۷۷	شکل ۳-۴ درخت Trie
۷۸	شکل ۴-۴ تابع Find_Candidate
۸۰	شکل ۵-۴ شکل ۵-۴ تابع ذخیره دیتاست در حافظه
۸۰	شکل ۶-۴ شبه کد الگوریتم سی ان دی آی، او دی آی ام

- شکل ۴-۷ شبه کد بازیافت مجموعه آیت‌های فراوان از CNDI ۸۱
- شکل ۵-۱ نمایش فشردگی برای پایگاه داده های تراکنشی به صورت باینری ۸۴
- شکل ۵-۲ مقایسه زمانی بین CNDIODAM و ODAM با تعداد ندهای متفاوت ۸۶
- شکل ۵-۳ مقایسه زمانی بین CNDIODAM و ODAM با مینیمم آستانه پشتیبانی متفاوت ۸۷
- شکل ۵-۴ مقایسه تعداد پیامهای انتقالی بین CNDIODAM و الگوریتم ODAM با پنج ند ۸۹
- شکل ۵-۵ حجم پیامهای انتقالی بین سایتها در دیتاست connect با آستانه فراوانی ۷۵٪ ۸۹

## فهرست جداول

صفحه	عنوان
۷	جدول ۱-۱ یک پایگاه داده نمونه
۱۰	جدول ۲-۱ نماد های مهم مربوط به الگوریتم ای. پریوری
۲۹	جدول ۱-۲ پایگاه داده D
۳۶	جدول ۲-۲ پایگاه داده تراکنشی
۴۸	جدول ۳-۲ پایگاه داده تراکنشی
۸۴	جدول ۱-۵ خصوصیات دیتاست ها

## فصل اول

### داده کاوی و کاوش قواعد وابستگی

## ۱- داده کاوی و کاوش قواعد وابستگی

### ۱-۱ مقدمه:

امروزه حجم زیادی از اطلاعات در پایگاه های داده مربوط به شرکت ها، مراکز تجاری و دولتی ذخیره می شود. استفاده معمول از این داده ها انجام عملیات گزارش گیری برای کاربران و مدیران است. استفاده دیگری که امروزه از حجم انبوه داده های ذخیره شده در پایگاه های داده و انبار های داده می شود انجام عملیات داده کاوی است. در عملیات داده کاوی ما به دنبال الگوهای پنهان و احتمالاً سودمند هستیم. برخی از این الگوها در انجام تصمیم گیری ها می توانند به مدیران کمک کرده و یا برای کاربران و مشتریان مفید باشند. الگوهایی که در عملیات مختلف داده کاوی پیدا می شوند انواع گوناگونی دارند. یک نوع پرکاربرد و معروف از این الگوها قوانین یا قواعد وابستگی<sup>۱</sup> هستند. معروف ترین کاربرد قوانین وابستگی در تحلیل سبد خرید برای فروشگاه ها و مراکز تجاری است. به عنوان مثال پس از کاوش پایگاه داده مربوط به یک فروشگاه زنجیره ای ممکن است مشخص شود که مشتریانی که از این فروشگاه شیر می خردند به احتمال ۶۰٪ کره نیز خواهند خرید. یافتن چنین قواعدی می تواند در چیدن قفسه ها، راهنمایی مشتریان و مسائل مدیریتی سودمند باشد. عملیات یافتن قواعد وابستگی را کشف یا کاوش قواعد وابستگی گویند.

امروزه استفاده از پایگاه داده های توزیع شده<sup>۲</sup> به ویژه برای سازمان ها و مراکزی که در یک گستره ی جغرافیایی قرار دارند، گسترش روز افزونی پیدا کرده است. با ذخیره شدن حجم زیادی از داده ها در این نوع پایگاه داده ها همانند پایگاه داده های متمرکز بحث داده کاوی در آنها مطرح شده است. کاوش قواعد وابستگی به عنوان یکی از پرکاربردترین روش های داده کاوی در پایگاه های داده توزیع شده از اهمیت خاصی برخوردار است. ما در این پایان نامه در ابتدا مروری بر چگونگی انجام عملیات داده کاوی خواهیم داشت. از میان عملیات مختلف داده کاوی کاوش قواعد وابستگی را به دقت بررسی خواهیم کرد. روش ها و الگوریتم های ارائه شده برای کاوش قواعد وابستگی در حالت متمرکز را بررسی خواهیم نمود. سپس الگوریتمهای کاوش قواعد وابستگی فشرده و در ادامه انجام عملیات داده کاوی در محیط توزیع شده را مورد مطالعه قرار داده و تفاوت های آن با حالت متمرکز را بررسی می کنیم. کاوش قواعد وابستگی در محیط توزیع شده به عنوان هدف اصلی این پایان نامه مدنظر قرار گرفته است و روش ها و الگوریتم های ارائه شده در این زمینه مورد مطالعه دقیق قرار گرفته و با هم مقایسه می شوند. در نهایت ما الگوریتم جدیدی برای

<sup>۱</sup> Association Rules Mining

<sup>۲</sup> Distributed DataBases

کاوش قوانین وابستگی در محیط توزیع شده ارائه خواهیم داد که از کارایی بهتری نسبت به الگوریتم های قبلی برخوردار باشد. الگوریتم جدید را پیاده سازی کرده و بر روی داده های توزیع شده ارزیابی خواهیم کرد.

## ۱-۲ داده کاوی

اخیراً سرعت تولید و جمع آوری داده ها در پایگاههای داده به صورت روزافزونی زیاد شده است. استفاده گسترده از بارکد در فروش تولیدات، کامپیوتری شدن تعداد زیادی از کارهای تجاری، اداری و دولتی و پیشرفت در زمینه ابزار جمع آوری داده ها ما را با حجم زیادی از داده ها مواجه کرده است. امروزه پایگاه داده ها در زمینه های تجاری، اداری، علمی، مهندسی و زمینه های دیگر استفاده می شوند. تعداد چنین پایگاههای داده ای به دلیل نیاز مبرم به جمع آوری و گزارش گیری از داده ها و همچنین وجود سیستم های پایگاه داده قدرتمند در حال افزایش است. این چنین رشد فزاینده ای در داده ها و پایگاههای داده یک نیاز ضروری برای ابزار جدیدی که بتوانند به طور هوشمند و خودکار این داده ها را پردازش کرده و به اطلاعات و به دانش های سودمند تبدیل کنند، بوجود آورده است. در نتیجه داده کاوی به یک زمینه تحقیقاتی با اهمیت فراوان تبدیل شده است.

از داده کاوی همچنین به عنوان کشف دانش در پایگاههای داده یاد می شود، که به معنی فرآیند استخراج اطلاعات غیر صریح و احتمالاً سودمندی از پایگاههای داده است که در گذشته نا شناخته و پنهان بوده اند. با انجام عملیات داده کاوی دانش های جالب و گاه غیر منتظره، نظم ها و الگوهای پنهان، یا اطلاعات سطح بالا می توانند از مجموعه ای از داده های مرتبط در پایگاه داده استخراج شوند و از زوایای مختلف مورد بررسی قرار گیرند. بنابراین پایگاههای داده حجیم را می توان به عنوان منبعی غنی و قابل اطمینان برای تولید و واریسی برخی دانش ها و اطلاعات در نظر گرفت.

کاوش اطلاعات و دانش از پایگاههای داده حجیم به عنوان یک موضوع کلیدی برای محققینی که در زمینه پایگاههای داده و یادگیری ماشین کار می کنند و به فرصتی برای کسب درآمد های بیشتر توسط شرکت های صنعتی و تجاری تبدیل شده است. دانش های کشف شده توسط داده کاوی می توانند در مدیریت اطلاعات، پردازش گزارش ها، انجام تصمیم گیری ها و بسیاری زمینه های دیگر استفاده شوند. به علت وجود گسترده داده ها در حجم زیاد و نیاز مبرم به تبدیل این داده ها به اطلاعات و دانش مفید برای کاربرد های مختلف، داده کاوی در سال های اخیر توجه زیادی را به خود جلب کرده است [1,2].

داده کاوی موضوعی وابسته به کاربرد است و کاربردهای مختلف نیازمند روش ها و تکنیک های داده کاوی مختلفی هستند. کاوش قواعد وابستگی، دسته بندی<sup>۳</sup>، خوشه بندی<sup>۴</sup>، پیش بینی<sup>۵</sup> و تحلیل سری های زمانی<sup>۶</sup> از جمله مهمترین روش ها و تکنیک های داده کاوی به شمار می آیند. در ادامه هر کدام از این روش ها به صورت خلاصه توضیح داده می شوند.

کشف یا کاوش قواعد وابستگی در پایگاه داده های رابطه ای<sup>۷</sup> یا تراکنشی<sup>۸</sup> که موضوع اصلی این مطالعه نیز است، اخیراً جذابیت زیادی را در انجمن های مربوط به پایگاههای داده بوجود آورده است. در این تکنیک داده کاوی وابستگی ها و ارتباطات بین داده های موجود در یک پایگاه داده بدست می آیند. نتیجه این عملیات داده کاوی دسته ای از قواعد است که به آنها قواعد وابستگی گفته می شود.

یکی دیگر از روش های مهم داده کاوی توانایی انجام عملیات دسته بندی در حجم زیاد داده هاست. این عملیات کاوش قوانین دسته بندی نیز نامیده می شود. در این روش اشیاء موجود در یک پایگاه داده بر اساس مقادیر چند خصوصیت از آنها، به دسته های مجزا تقسیم می شوند. در دسته بندی داده ها مجموعه ای از داده های آزمایشی تحلیل می شوند. برای مجموعه داده های آزمایشی برچسب کلاس ها مشخص است. برای هر کلاس داده های آزمایشی مدلی بر اساس خصوصیات داده ها ساخته می شود. حاصل عملیات دسته بندی می تواند یک درخت تصمیم یا مجموعه ای از قوانین دسته بندی باشد که برای فهم بهتر داده های موجود در پایگاه داده و همچنین دسته بندی داده هایی که در آینده به پایگاه داده اضافه می شوند به کار می رود. دسته بندی داده ها ارتباط تنگاتنگی با کاوش قواعد وابستگی دارد، به طوری که گاهی دسته بندی را به کمک قواعد وابستگی انجام می دهند.

به عنوان مثال برای فروشندهی ماشین دسته بندی مشتریان بر اساس تمایل و علاقه آنها به انواع مختلف ماشین مطلوبست به طوری که بتواند به مشتریان خدمات بهتری ارائه دهد و کاتالوگ های محصولات جدید مورد نظر آنها را برایشان بفرستد و درآمد حاصل از فروش خود را افزایش دهد.

خوشه بندی داده ها از مهمترین روش های داده کاوی به شمار می آید. درخوشه بندی مجموعه ای از داده ها گروهبندی می شوند. فرق خوشه بندی با دسته بندی داده ها در این است

<sup>3</sup> Classification

<sup>4</sup> Clustering

<sup>5</sup> Prediction

<sup>6</sup> Time Series Analysis

<sup>7</sup> Relational

<sup>8</sup> Transactional



که در خوشه بندی برخلاف دسته بندی تعداد کلاس ها در ابتدا مشخص نیستند. خوش بندی داده ها بر اساس اصل مفهومی زیر صورت می گیرد:

"حداکثر کردن شباهت های بین اعضای هر کلاس و حداقل کردن شباهت ها بین اعضای مربوط به کلاس های مختلف".

به عنوان مثالی از خوشه بندی، مجموعه ای از کالاها را می توان در ابتدا به صورت مجموعه ای از کلاس های مختلف خوشه بندی کرده و سپس مجموعه ای از قوانین را بر اساس این چنین دسته بندی نتیجه گیری کرد. پیش بینی یکی دیگر از تکنیک های داده کاوی محسوب می شود که در آن مقادیر ممکن برای متغیرهای نامعلوم پیش بینی می شوند. در پیش بینی ابتدا داده هایی که به متغیر نامعلوم مربوط هستند بوسیله برخی تحلیل های آماری پیدا می شوند. سپس از برخی روش های هوشمند مانند شبکه های عصبی و الگوریتم ژنتیک برای انجام پیش بینی استفاده می شود. برای مثال مقدار حقوقی که یک کارمند می گیرد را می توان با استفاده از چگونگی توزیع حقوق کارمندان مشابه در آن اداره، پیش بینی کرد. روش های دیگری از جمله تحلیل رگرسیون<sup>9</sup>، تحلیل وابستگی<sup>10</sup>، درخت تصمیم<sup>11</sup> در انجام یک پیش بینی با کیفیت مؤثرند. تحلیل سری های زمانی از دیگر روش های کاربردی داده کاوی است. در این روش حجم زیادی از داده های سری زمانی برای یافتن خصوصیات جالب توجه و نظم های مشخص، تحلیل می شوند. رخداد وقایع متوالی، مجموعه وقایعی که بعد از یک واقعه مشخص به وقوع می پیوندند، روند ها و انحراف ها از جمله این نظم ها و پدیده های جالب توجه هستند. برای مثال می توان روند تغییر قیمت برای یک کالای مشخص را در یک کارخانه با استفاده از داده های گذشته، شرایط تجاری و رقبای تجاری، پیش بینی کرد.

موارد ذکر شده در بالا تعدادی از تکنیک ها و روش های مهم در داده کاوی بودند روش های دیگری نیز وجود دارند که در این نوشته مجالی برای بحث در مورد آنها نیست.

### ۱-۳ مراحل داده کاوی

به داده کاوی عملیات اکتشاف دانش در پایگاه های داده هم گفته می شود، اگرچه برخی محققین داده کاوی را بخشی از عملیات کلی تر اکتشاف دانش<sup>12</sup> می دانند. به طور کلی عملیات اکتشاف دانش شامل اجرای تکراری مراحل زیر است:

<sup>9</sup> Regression Analysis

<sup>10</sup> Correlation Analysis

<sup>11</sup> Decision Tree

<sup>12</sup> Knowledge Discovery

- خالص کردن داده ها<sup>13</sup>: در این مرحله داده ها از وجود خطاها، داده های مشکوک و داده های غیر مرتبط پاک می شوند.
  - سرهم سازی داده ها<sup>14</sup>: در این مرحله داده هایی که از منابع مختلف گردآوری شده اند در یک منبع واحد مجتمع می شوند.
  - گزینش داده ها<sup>15</sup>: در این مرحله داده هایی را که عملیات تحلیل روی آنها انجام می شود از پایگاه داده بازیابی می شوند.
  - تبدیل داده ها<sup>16</sup>: در این مرحله داده ها به قالب و فرم مناسب برای انجام عملیات داده کاوی تبدیل می شوند.
  - داده کاوی: مرحله ای ضروری است که در آن روش های هوشمند به منظور استخراج الگوها، به داده های تبدیل یافته اعمال می شوند.
  - ارزیابی الگوها<sup>17</sup>: در این مرحله الگوهای بدست آمده توسط یک سری معیارهای مقبولیت سنجیده می شوند و در نهایت از میان الگوهای بدست آمده آنهایی که واقعا با ارزش هستند تشخیص داده می شوند.
  - ارائه دانش<sup>18</sup>: در این مرحله نتایج بدست آمده به شکل معنا داری به وسیله ی عملیات گرافیکی و سایر روش های نمایش دانش به کاربر عرضه می گردند.
- با گسترش روز افزون استفاده از پایگاه های داده و انبار های داده چهار مرحله ی نخست شامل خالص کردن، سرهم سازی، گزینش و تبدیل، با ساختن انبارهای داده و انجام برخی عملیات گزارش گیری مربوط به انبارهای داده قابل انجام هستند. عملیات داده کاوی، ارزیابی الگوها و ارائه دانش گاهی در یک مرحله کلی بنام داده کاوی انجام می گیرد.

#### ۴-۱ روشهای داده کاوی

داده، انبارهای داده، آمار، یادگیری ماشین، نمایش داده، بازیابی اطلاعات و محاسبات سریع سرچشمه گرفته است، انجام عملیات مختلف داده کاوی نیازمند آگاهی از این زمینه هاست. آگاهی از فنون مختلف هوش مصنوعی مانند شبکه های عصبی، محاسبات فازی و الگوریتم ژنتیک برای انجام برخی عملیات داده کاوی ضروری است.

<sup>13</sup> Data Cleaning

<sup>14</sup> Data integration

<sup>15</sup> Data selection

<sup>16</sup> Data transformation

<sup>17</sup> Pattern evaluation

<sup>18</sup> Knowledge presentation

مجموعه وسیعی از عملیات مربوط به آمار که در طول سال ها توسعه یافته اند در انجام برخی عملیات تحلیل داده ها استفاده می شوند. یادگیری ماشین در عملیاتی چون دسته بندی داده ها استفاده می شود. شبکه های عصبی در عملیات دسته بندی، خوشه بندی و پیش بینی کاربرد فراوانی دارد. به هر حال از آنجایی که در عملیات داده کاوی ما با حجم زیاد داده هایی که در پایگاه های داده ذخیره می شوند سر و کار داریم در استفاده از این روش ها به مشکل کارائی و مقیاس پذیری روبرو می شویم.

استفاده از ساختمان داده های مناسب، روش های اندیس گذاری و دستیابی به داده ها که در بحث پایگاه های داده مطرح شده اند تأثیر بسزایی در انجام سریع عملیات داده کاوی دارند. روش های تحلیل داده ها که در آمار و یادگیری ماشین مطرح هستند برای روبرو شدن با حجم زیاد داده ها در داده کاوی، بازنگری می شوند. برای انجام داده کاوی مؤثر باید الگوریتم های مقیاس پذیر و مبتنی بر مجموعه ها داشته باشیم.

برخلاف عملیات تحلیل داده ها که همیشه یک فرض وجود دارد، در عملیات داده کاوی هدف ما اکتشاف است و سعی می کنیم الگوهایی ناشناخته از داده ها را بدست آوریم و نیازمند جستجوهای سنگین و وقت گیر در داده ها هستیم. بنابراین انجام محاسبات سریع نقش مهمی را در داده کاوی بازی می کند. روش های موازی، توزیع شده و افزایشی بیشتری در زمینه داده کاوی به این منظور باید توسعه یابند.

## ۱-۵ کاوش قواعد وابستگی

یکی از مهمترین و پرکاربردترین روش های داده کاوی قواعد وابستگی است، که در آن هدف پیدا کردن قوانینی است که در پایگاه داده های با حجم بالا وجود دارند. به عنوان مثال پایگاه داده مربوط به یک فروشگاه زنجیره ای کالا را در نظر بگیرید که در آن اطلاعات مربوط به اجناس خریداری شده توسط هر مشتری ذخیره می شود، در طول سالهای کار این فروشگاه حجم داده های ذخیره شده در پایگاه داده انباشته و انباشته تر می شود. استفاده معمول از این داده ها گزارش گیری های مربوط به وضعیت فروش و مشتریان این فروشگاه است. اما استفاده مفید دیگری که از این داده ها می توان کرد یافتن یک سری قوانین است که در آن وابستگی و ارتباط بین کالای خریداری شده توسط مشتریان نمایان می شود. به عنوان مثال جدول ۱-۱ را به عنوان پایگاه داده یک فروشگاه در نظر بگیرید، ستون سمت چپ شماره تراکنش هایی است که در این پایگاه داده ذخیره شده اند، و ستون ستون سمت راست معرف مجموعه اقلام مربوط به هر تراکنش است.

جدول ۱-۱ یک پایگاه داده نمونه

Tid	Items
1	{bread, coke, milk}
2	{bread, butter, milk, ice cream}
3	{ice cream, coke}
4	{battery, bread, butter, milk}
5	{bread, butter, milk}
6	{battery, ice cream, bread, butter}

اگر چه در پایگاه داده مثال ما تنها شش تراکنش وجود دارد، ولی در پایگاههایی که روی آنها عملیات کاوش قواعد وابستگی صورت می گیرد حجم زیادی از تراکنش ها وجود دارد. یکی از این قوانین موجود در پایگاه داده جدول ۱-۱ می تواند به این صورت باشد که هرکسی که از این فروشگاه کره خریده است به احتمال ۸۰٪ نان هم خریده است، که این قاعده در ۶۶٪ موارد خرید یا تراکنش های این فروشگاه دیده شده است. این نوع قوانین را در اصطلاح قواعد وابستگی می نامند و مجموعه اقلامی که به طور مکرر در پایگاه داده تکرار می شوند و ما قوانین مورد نظر را از روی آنها تولید می کنیم مجموعه اقلام تکراری<sup>۱۹</sup> یا مجموعه الگوهای مکرر<sup>۲۰</sup> نام دارند مانند مجموعه {bread, butter} در پایگاه داده مثال بالا. احتمال مطرح شده را درصد اطمینان<sup>۲۱</sup> و درصدی از تراکنش ها که این قاعده را پوشش می دهند، درصد پشتیبانی<sup>۲۲</sup> این قاعده می نامند. عملیات مربوط به پیدا کردن این قوانین را کاوش قواعد وابستگی یا به اختصار ARM می نامند [3]. درصد پشتیبانی و درصد اطمینان، دو معیار مربوط به مقبولیت قوانین کشف شده هستند. معیار های دیگری نیز در این زمینه وجود دارند.

معمولاً برای یافتن قوانین حداقلی برای اندازه معیار های پشتیبانی و اطمینان توسط کاربر در نظر گرفته می شود، که به آنها به ترتیب حداقل حد آستانه<sup>۲۳</sup> پشتیبانی و حداقل حد آستانه اطمینان گفته می شود. عملیات اصلی در مورد کاوش قواعد وابستگی یافتن مجموعه اقلام تکراری با استفاده از حداقل حد آستانه پشتیبانی است. در واقع با داشتن این مجموعه اقلام تولید قواعد وابستگی بسیار سر راست خواهد بود. بنابراین مسئله کاوش قواعد وابستگی به مسئله یافتن مجموعه اقلام تکراری کاهش می یابد. از آنجایی که کاوش مجموعه اقلام تکراری نیاز به خواندن از

<sup>19</sup> Frequent itemset

<sup>20</sup> Frequent pattern

<sup>21</sup> Confidence

<sup>22</sup> Support

<sup>23</sup> Threshold