

الله
يَا مُحَمَّدُ



پرديس بين الملل ارس

گروه علوم کامپیووتر

پایان نامه

برای دریافت درجه کارشناسی ارشد در رشته علوم کامپیووتر

عنوان

کشف قوانین انجمنی به صورت موازی بر اساس یک روش خوشبندی جدید
و با استفاده از یک الگوریتم تکاملی

استاد راهنما

دکتر شهریار لطفی

استاد مشاور

دکتر جابر کریمپور ینگجه

پژوهشگر

لیلا امیری

تابستان ۱۳۹۳

با تشکر و قدردانی از

استاد بزرگوار جناب آقای دکتر شهریار لطفی که با زحمات بیدریغ و دلسوزانه خویش راهنمایی پایان نامه را بر عهده داشتند و جناب آقای دکتر جابر کریمپور ینگجه که مشاوره‌ی این پایان نامه را تقبل فرمودند. همچنین پدر و مادر دلسوزم که مرا در پیشبرد پایان نامه یاری نمودند، تشکر و قدردانی می‌نمایم.

لیلا امیری

شهریور ۱۳۹۳

نام خانوادگی: امیری	نام: لیلا
عنوان پایان نامه: کشف قوانین انجمنی به صورت موازی بر اساس یک روش خوشبندی جدید و با استفاده از یک الگوریتم تکاملی	
استاد راهنما: دکتر شهریار لطفی	
استاد مشاور: دکتر جابر کریم پور ینگجه	
مقطع تحصیلی: کارشناسی ارشد رشته: علوم کامپیوتر گرایش: سیستم‌های کامپیوترا	
دانشگاه: تبریز دانشگاه: پردیس بین المللی ارس تاریخ فارغ التحصیلی: ۱۳۹۳ تعداد صفحات: ۱۳۷	
کلید واژه‌ها: کشف قوانین انجمنی، خوشبندی، الگوریتم تکاملی و الگوریتم رقابت استعماری	
چکیده:	
<p>کشف قوانین انجمنی، یکی از روش‌های مهم و پرکاربرد داده‌کاوی جهت کشف دانش نهفته در داده‌ها است که با استفاده از آن می‌توان روابط و وابستگی‌های مفیدی که در مجموعه‌های بزرگی از اقلام داده موجود می‌باشند را کشف نمود. این روابط و وابستگی‌ها در تصمیم‌گیری‌ها، نقش مهمی دارند. یافتن چنین روابطی، داخل مجموعه داده‌ها به دلیل ماهیت نمایی آن و حجم بالای داده‌ها بسیار زمان بر است. در این پایان نامه برای کشف قوانین انجمنی، ابتدا با ارائه یک روش خوشبندی مبتنی بر الگوریتم رقابت استعماری، تراکنش‌ها به خوشبندی مناسب تقسیم می‌شوند. سپس یک روش تکاملی بر پایه الگوریتم رقابت استعماری برای کشف قوانین انجمنی ارائه می‌گردد که این الگوریتم به طور جداگانه و مستقل بر روی هر یک از خوشبندی‌ها اجرا می‌شود. در نهایت، قوانین به دست آمده از همه خوشبندی‌ها در یک جا جمع‌آوری شده و قوانین نهایی تولید می‌گرددند. نتایج حاصل از آزمایشات مختلف بر روی چند مجموعه داده شناخته شده، کارآیی راهکار پیشنهادی را در کشف قوانین با دقت مناسب و کاهش هزینه‌ها تایید می‌کند.</p>	

فهرست مطالب

صفحه	عنوان
	فصل ۱ - مقدمه ۱
۵	فصل ۲ - شرح مسئله ۵
۶	۱-۱ بیان مسئله ۶
۸	۲-۱ ورودی و خروجی ۸
۸	۳-۱ فرضهای مسئله ۸
۸	۴-۱ هدف ۸
۸	۵-۱ نوآوری ۸
۹	۶-۱ سوالات تحقیق ۹
۹	۷-۱ فرضیات تحقیق ۹
۹	۸-۱ خلاصه فصل ۹
۱۱	فصل ۳ - مفاهیم پایه‌ای ۱۱
۱۱	۱-۱ داده کاوی ۱۱
۱۳	۲-۱ طبقه‌بندی ۱۳
۱۳	۳-۱ خوشبندی ۱۳
۱۴	۱-۲-۱ اندازه‌گیری شباهت ۱۴
۱۷	۲-۲-۱ طبقه‌بندی الگوریتم‌های خوشبندی ۱۷
۲۱	۳-۱-۱ K-Means ۲۱
۲۲	۴-۱-۱ الگوریتم K-Medoids ۲۲
۲۲	۵-۱-۱ الگوریتم K-Modes ۲۲
۲۳	۶-۱-۱ اعتبار خوش ۲۳
۲۵	۴-۱-۱ کشف قوانین انجمنی ۲۵
۲۵	۴-۲-۱ اصطلاحات و تعاریف ۲۵

صفحه	عنوان
۲۶	۴-۳ قوانین انجمنی ۲-۴
۲۹	۴-۳ انواع الگوریتم‌های کشف قوانین انجمنی ۶-۴
۳۱	۳-۶ بهینه‌سازی ۶-۶
۳۲	۳-۶ بهینه‌سازی چند هدفه ۱-۶
۳۳	۳-۶ بررسی روش‌های بهینه‌سازی ۲-۶
۳۹	۳-۶ نسخه گسسته الگوریتم رقابت استعماری ۴-۶
۴۰	۳-۷ خلاصه فصل ۷
۴۱	۴-۸ فصل ۴ - راهکارهای گذشته
۴۲	۴-۸ روش‌های سنتی ۱-۴
۴۲	۴-۸ ۱-۱ الگوریتم AIS
۴۳	۴-۸ ۲-۱ الگوریتم Apriori
۴۴	۴-۸ ۲-۲ روش‌های بهبود فرآیند کشف قوانین انجمنی
۵۰	۴-۸ ۳-۳ شیوه‌ای مبتنی بر الگوریتم ICA جهت استخراج قوانین انجمنی
۵۳	۴-۸ ۴-۳ روش‌های کشف قوانین انجمنی مبتنی بر بخش‌بندی با استفاده از خوش‌بندی
۵۳	۴-۸ ۴-۱ کشف قوانین انجمنی نادر با استفاده از خوش‌بندی بر اساس Seedها
۵۵	۴-۸ ۴-۲ کشف قوانین با استفاده از الگوریتم ژنتیک بر اساس خوش‌بندی
۵۹	۴-۸ ۴-۳ کشف قوانین انجمنی فازی مبتنی بر خوش‌بندی
۶۸	۴-۸ ۵-۶ خلاصه
۶۸	۵-۶ فصل ۵ - راهکار پیشنهادی
۶۹	۵-۶ ۱-۱ مراحل راهکار پیشنهادی
۶۹	۵-۶ ۱-۱ ۱-۱ الگوریتم پیشنهادی برای خوش‌بندی تراکنش‌ها
۸۱	۵-۶ ۱-۱ ۲-۱ الگوریتم پیشنهادی برای کشف قوانین انجمنی

عنوان		صفحه
۳-۱-۵ جمع‌آوری قوانین تولید شده از کل خوشها و تولید قوانین نهایی ۹۲		۹۲
۲-۵ فلوچارت راهکار پیشنهادی ۹۲		۹۲
۳-۵ خلاصه فصل ۹۳		۹۳
فصل ۶ - ارزیابی و نتایج عملی ۹۵		۹۵
۱-۶ آزمایشات ۹۶		۹۶
۱-۱-۶ قابلیت اطمینان ۹۶		۹۶
۲-۱-۶ بررسی همگرایی ۹۹		۹۹
۳-۱-۶ بررسی پایداری ۱۰۲		۱۰۲
۲-۶ بررسی نتایج حاصل از راهکار پیشنهادی و مقایسه آن با روش‌های موجود ۱۰۴		۱۰۴
۱-۲-۶ بررسی نتایج الگوریتم پیشنهادی و مقایسه با الگوریتم Apriori ۱۰۴		۱۰۴
۲-۲-۶ مقایسه الگوریتم پیشنهادی و الگوریتم AICluster ۱۱۳		۱۱۳
۳-۲-۶ مقایسه الگوریتم پیشنهادی و الگوریتم MINICA ۱۱۵		۱۱۵
۳-۶ بحث ۱۱۶		۱۱۶
۴-۶ خلاصه فصل ۱۱۷		۱۱۷
۱-۷ نتیجه‌گیری و کارهای آتی ۱۱۸		۱۱۸
۲-۷ راهکارهای آتی ۱۲۱		۱۲۱
مراجع ۱۲۳		۱۲۳

فهرست جدول‌ها

عنوان	صفحه
جدول ۱-۲ نمونه‌ای از یک تراکنش [۲]	۷
جدول ۱-۳ توابع فاصله میان داده‌های x و y [۲]	۱۴
جدول ۲-۳ ترکیب‌های ممکن مقادیر x و y برای اندازه‌گیری شباهت [۲]	۱۶
جدول ۱-۴ محاسبه‌ی سرعت (بر اساس ثانیه) بر اساس تعداد رکوردها [۱۹]	۵۸
جدول ۲-۴ نتایج خوشبندی پس از اجرای الگوریتم خوشبندی [۶]	۶۴
جدول ۱-۶ پارامترهای راه‌کار پیشنهادی برای اجرای بر روی TestData	۹۸
جدول ۲-۶ پارامترهای الگوریتم پیشنهادی برای خوشبندی بر روی داده	۹۹
جدول ۳-۶ مقدار پارامترهای الگوریتم در مجموعه داده Zoo	۱۰۷
جدول ۴-۶ مقایسه نتایج اجرا الگوریتم‌های مختلف بر روی مجموعه داده Zoo	۱۰۷
جدول ۵-۶ مقدار پارامترها در مجموعه داده Data5000	۱۰۸
جدول ۶-۶ مقایسه نتایج حاصل از اجرای الگوریتم در Data5000	۱۰۹
جدول ۷-۶ پارامترهای در نظر گرفته شده برای مجموعه داده Mushroom	۱۱۱
جدول ۸-۶ پارامترهای در نظر گرفته شده برای مجموعه داده ChessData	۱۱۱
جدول ۹-۶ مقایسه نتایج حاصل از اجرای الگوریتم بر روی مجموعه داده Mushroom	۱۱۲
جدول ۱۰-۶ مقایسه نتایج حاصل از اجرای الگوریتم بر روی مجموعه داده ChessData	۱۱۲
جدول ۱۱-۶ مقدار پارامترهای الگوریتم در مجموعه داده Zoo	۱۱۳
جدول ۱۲-۶ نتایج مقایسه الگوریتم پیشنهادی و الگوریتم Apriori Inverse و AICluster	۱۱۴
جدول ۱۳-۶ مقدار پارامترهای الگوریتم در مجموعه داده Car Evaluation	۱۱۵
جدول ۱۴-۶ نتایج مقایسه الگوریتم پیشنهادی و الگوریتم MINICA	۱۱۵

فهرست شکل‌ها

صفحه	عنوان
۱۵	شکل ۱-۳ داده‌های دودویی x و y [۲]
۲۷	شکل ۲-۳ مثالی از یک سلسله مراتب برای مواد غذایی [۲]
۳۶	شکل ۳-۳ شمای کلی فرآیند سیاست جذب
۳۷	شکل ۴-۳ شمای کلی فرآیند جابه‌جایی کشور استعمارگر و مستعمره
۳۸	شکل ۵-۳ شمای کلی فرآیند رقابت استعماری میان امپراطوری‌ها
۳۹	شکل ۶-۳ شمای کلی فرآیند سقوط امپراطوری‌های ضعیف
۵۱	شکل ۱-۴ نحوه کدگذاری قانون در یک کشور [۸]
۵۱	شکل ۲-۴ نحوه نمایش قانون در یک کشور [۸]
۵۷	شکل ۳-۴ ساختار کلی روش پیشنهادی [۱۹]
۵۸	شکل ۴-۴ کدگذاری کروموزوم برای قانون [۱۹] $ACF \rightarrow BE$
۷۱	شکل ۱-۵ نحوه نمایش یک کشور
۷۱	شکل ۲-۵ نحوه نمایش یک کشور
۷۴	شکل ۳-۵ A، کشور استعمارگر و B، کشور مستعمره
۷۴	شکل ۴-۵ نتیجه محاسبه A-B
۷۴	شکل ۵-۵ نتیجه محاسبه B-A
۷۵	شکل ۶-۵ Add و Del پس حذف برخی از عناصر
۷۵	شکل ۷-۵ position جدید برای B پس از اعمال سیاست جذب
۷۶	شکل ۸-۵ عناصر انتخاب شده از Position کشور برای تغییر
۷۶	شکل ۹-۵ Position جدید کشور پس از اجرای انقلاب با استفاده از روش اول
۷۷	شکل ۱۰-۵ Position جدید کشور پس از اجرای انقلاب با استفاده از روش دوم
۸۲	شکل ۱۱-۵ نمایش قانون $AD \rightarrow BCE$

عنوان

صفحه

..... شکل ۱۲-۵ نحوه نمایش قانون $AB \rightarrow CD$	۸۴
..... شکل ۱۳-۵ نحوه نمایش قانون $E \rightarrow AC$	۸۴
..... شکل ۱۴-۵ مجموعه‌های Addl و Removel	۸۵
..... شکل ۱۵-۵ Addl و Removel پس حذف برخی از عناصر	۸۵
..... شکل ۱۶-۵ مجموعه‌های AddR و RemoveR	۸۵
..... شکل ۱۷-۵ AddR و RemoveR پس از حذف برخی از عناصر	۸۶
..... شکل ۱۸-۵ کشور مستعمره پس از اعمال سیاست جذب	۸۶
..... شکل ۱۹-۵ مقداردهی اولیه آرایه Chrom	۸۸
..... شکل ۲۰-۵ قرار گرفتن اقلام موجود در قانون $AB \rightarrow C$ به طور تصادفی در آرایه chrom	۸۸
..... شکل ۲۱-۵ عناصر انتخاب شده از آرایه chrom برای تغییر	۸۸
..... شکل ۲۲-۵ تغییر مکان‌های انتخابی با محتوای chrom در آرایه	۸۹
..... شکل ۲۳-۵ تغییر مکان‌های انتخابی با محتوای غیر صفر در آرایه chrom	۸۹
..... شکل ۲۴-۵ نحوه نمایش قانون تولید شده پس از اجرای عمل گر انقلاب بر روی قانون	۹۰
..... شکل ۲۵-۵ فلوچارت راه کار پیشنهادی	۹۳
..... شکل ۱-۶ نتایج اجرای الگوریتم Apriori بر روی TestData	۹۷
..... شکل ۲-۶ نتایج اجرای الگوریتم پیشنهادی بر روی TestData	۹۸
..... شکل ۳-۶ نمودار هم‌گرایی الگوریتم پیشنهادی برای خوشبندی بر روی داده TestData	۱۰۰
..... شکل ۴-۶ نمودار هم‌گرایی الگوریتم کشف قوانین انجمنی بر روی داده TestData	۱۰۲
..... شکل ۵-۶ نمودار پایداری الگوریتم پیشنهادی برای خوشبندی بر روی داده TestData	۱۰۳
..... شکل ۶-۶ نتایج آزمون پایداری الگوریتم پیشنهادی برای خوشبندی تراکنش‌ها	۱۰۳
..... شکل ۷-۶ نمودار پایداری الگوریتم پیشنهادی کشف قوانین انجمنی بر روی داده TestData	۱۰۴

عنوان

صفحه

- شكل ٦-٨ نتایج آزمون پایداری الگوریتم کشف قوانین انجمنی بر روی داده TestData ١٠٤
- شكل ٩-٦ نتایج الگوریتم Apriori بر روی مجموعه داده Zoo ١٠٥
- شكل ١٠-٦ نتایج الگوریتم پیشنهادی بر روی مجموعه داده Zoo ١٠٦
- شكل ١١-٦ نتایج الگوریتم پیشنهادی بدون مرحله پسپردازش بر روی مجموعه داده Zoo ١٠٦
- شكل ١٢-٦ نتایج اجرای الگوریتم پیشنهادی غیرخوشبندی بر روی مجموعه داده Zoo ١٠٦
- شكل ١٣-٦ نتایج الگوریتم بدون مرحله خوشبندی بر روی مجموعه داده Data5000 ١٠٨
- شكل ١٤-٦ نتایج الگوریتم پیشنهادی بر روی مجموعه داده Data5000 ١٠٩
- شكل ١٥-٦ نتایج الگوریتم پیشنهادی بر روی مجموعه داده Zoo ١١٣
- شكل ١٦-٦ نمودار همگرایی الگوریتم پیشنهادی و الگوریتم MINICA ١١٦

فصل ا

مقدمہ

در سال‌های گذشته، پیشرفت‌های قابل توجهی در روش‌های جمع‌آوری، ذخیره‌سازی و انتقال حجم عظیمی از داده‌ها در زمینه‌های مختلف، صورت گرفته است. با توجه به افزایش چشم‌گیر حجم داده‌ها، نیاز به روش‌های بهتر، سریع‌تر و ارزان‌تر جهت استخراج اطلاعات و دانش نهفته‌ی موجود در این حجم از داده‌ها می‌باشد، و گرنه داده‌های موجود، فاقد ارزش خواهند بود. داده‌کاوی بخشی از فرآیند کشف دانش از داده‌ها می‌باشد. داده‌کاوی به معنای یافتن نیمه‌خودکار الگوهای پنهان موجود در مجموعه‌ی داده‌های موجود می‌باشد [۱]. از طرف دیگر صاحبان داده، درک کمی از داده‌ها دارند و دانش آن‌ها نسبت به این داده‌ها کم است. یکی از مسائل اساسی در داده‌کاوی این است که با وجود حجم بالای داده‌ها، یک مدل کوچک و نه زیاد پیچیده از داده‌ها ایجاد شود که در عین حال که داده‌ها را به خوبی توصیف نماید، همچنین ساده و قابل فهم باشد. داده‌کاوی در حوزه‌های تصمیم‌گیری، پیش‌بینی و تخمین مورد استفاده قرار می‌گیرد. کشف قوانین انجمنی یکی از روش‌های هدایت نشده و بسیار مهم در داده‌کاوی می‌باشد. با توجه به افزایش روز افزون حجم داده‌ها و مقیاس‌پذیری الگوریتم‌های کشف قوانین انجمنی، این روش یکی از پرکاربردترین شیوه‌های داده‌کاوی است. قوانین انجمنی در زمینه‌های مختلفی از جمله تجارت (مانند بازاریابی و تعیین راهبرد برای قیمت‌گذاری)، پیش‌پردازش داده‌ها، خوشبندی، سیستم‌های شخصی‌سازی و پیشنهاد برای جستجوی صفحات وب (مانند پیشنهادهایی که سایت آمازون برای خرید کتاب‌های مرتبط ارائه می‌دهد) کاربرد دارد [۲].

هدف اصلی از کشف قوانین انجمنی، یافتن روابط بین داده‌ها با استفاده از تحلیل داده‌هایی است که در بیشتر موارد به طور هدایت نشده جمع‌آوری شده‌اند. این روابط و وابستگی‌ها در تصمیم‌گیری‌ها نقش مهمی دارند. در روش‌های سنتی، کشف قوانین انجمنی از دو مرحله اصلی

تشکیل شده است. ابتدا تمامی مجموعه‌های اقلام پر تکرار، تولید شده و در گام بعدی قوانین انجمانی قوی که محدودیت‌های معینی را ارضا می‌کنند، بر اساس این مجموعه‌ها کشف می‌شوند. مرحله دوم دارای هزینه‌ی محاسباتی کمی می‌باشد اما مرحله اول دارای هزینه‌ی محاسباتی نمایی است. از طرفی در تولید قوانین انجمانی، داده‌ها ممکن است چندین بار پیمایش شوند. حجم داده‌ها در مسائل مربوط به داده‌کاوی بالا می‌باشد و در بیشتر موارد تمامی داده‌ها، قابل بارگذاری در حافظه‌ی اصلی نمی‌باشند بنابراین سربار ناشی از عملیات ورودی و خروجی در این روش‌ها بسیار بالا است .[۴] [۳]

روش‌های مختلفی برای کاهش هزینه‌های مذکور و بهبود فرآیند کشف قوانین انجمانی ارائه شده است. در سال‌های گذشته، استفاده از الگوریتم‌های تکاملی در کشف قوانین انجمانی به دلیل تولید قوانین، بدون نیاز به تولید مجموعه‌های اقلام پر تکرار، بسیار مورد توجه قرار گرفته‌اند. هم‌چنان استفاده از روش‌های مبتنی بر خوشبندی تراکنش‌ها، به دلیل کاهش تعداد پیمایش داده‌ها و امکان اجرای موازی الگوریتم‌های کشف قوانین انجمانی در هر خوش، در کاهش هزینه‌ها و افزایش کارآیی الگوریتم‌های کشف قوانین انجمانی بسیار موثر می‌باشند.

در این پایان‌نامه برای کشف قوانین انجمانی، ابتدا با ارائه یک الگوریتم تکاملی جدید برای کشف قوانین انجمانی ارائه می‌گردد که این الگوریتم به طور جداگانه و مستقل بر روی هر یک از خوش‌ها اجرا می‌شود. در نهایت، قوانین به دست آمده از همه خوش‌ها در یک جا جمع‌آوری شده و قوانین نهایی تولید می‌گردند. راه کار پیشنهادی به اختصار ICA-ARMC^۱ نام‌گذاری شده است.

در ادامه در فصل دوم، به بیان مسئله‌ی کشف قوانین انجمانی، ورودی و خروجی، و فرض‌های مسئله، اهداف خواهیم پرداخت. در فصل سوم، مفاهیم و اصطلاحات پایه‌ای، تعاریف و روش‌های مختلف مورد استفاده در کشف قوانین انجمانی شرح داده می‌شود. در فصل چهارم، روش‌های سنتی و روش‌های ارائه شده برای بهبود فرآیند کشف قوانین انجمانی را مورد بررسی قرار خواهیم داد. در

^۱ ICA Association Rules Mining Using Clustering

ادامه در فصل پنجم، به تشریح راهکار پیشنهادی و مراحل مختلف آن خواهیم پرداخت. فصل ششم، شامل آزمایش‌ها و بررسی نتایج راهکار پیشنهادی و مقایسه با روش‌های موجود می‌باشد و در آخر در فصل هفتم، نتیجه‌گیری و پیشنهادات برای راهکارهای آتی را خواهیم داشت.

فصل ۲

شرح مسئلہ

کشف قوانین انجمنی یکی از مهمترین روش‌های داده‌کاوی می‌باشد. با توجه به افزایش چشم‌گیر داده‌ها، کشف قوانین انجمنی به دلیل درک راحت‌تر و مقیاس‌پذیری الگوریتم‌های آن یکی از پرکاربردترین روش‌های کشف قوانین انجمنی است. در این فصل، به بیان مسئله‌ی کشف قوانین انجمنی، ورودی و خروجی، و فرض‌های مسئله و اهداف خواهیم پرداخت.

۱-۲ بیان مسئله

کشف قوانین انجمنی^۱، یکی از روش‌های بسیار مهم داده‌کاوی^۲ جهت کشف دانش نهفته در داده‌ها می‌باشد. با استفاده از این شیوه، می‌توان روابط و وابستگی‌های مفیدی که در مجموعه‌های بزرگ از اقلام داده^۳ موجود می‌باشد را کشف نمود. اقلام داده، در قالب تراکنش‌ها^۴ ذخیره می‌شوند. تراکنش‌ها را می‌توان به کمک یک فرآیند خارجی تولید نمود، یا از پایگاه‌های داده یا انباره‌های داده استخراج نمود. یکی از کاربردهای بارز قوانین انجمنی، تحلیل سبد خرید^۵ می‌باشد که سعی می‌شود با یافتن وابستگی و روابط موجود بین اجناس خریداری شده به وسیله مشتری‌ها، الگوهای خرید، شناسایی و تحلیل شوند. به عنوان مثال، مشتری‌ها، شیر و نان را با هم می‌خرند. این نوع قوانین در افزایش سود و رقابت‌های بازاریابی استفاده می‌شود.

در مثال تحلیل سبد خرید، هر یک از محصولات موجود در فروشگاه، متناظر با یک قلم داده می‌باشد. سبد خرید متعلق به هر یک از مشتری‌ها، نشان دهنده‌ی اقلام خریداری شده به وسیله‌ی آن مشتری است. هر سبد خرید مربوط به یک مشتری به عنوان یک تراکنش در مجموعه‌ی داده

^۱ Association Rules Mining (ARM)

^۲ Data Mining (DM)

^۳ itemset

^۴ transactions

^۵ market basket analyse

می‌باشد و شامل یک شناسه‌ی یکتا است. که در جدول ۱-۲ نمونه‌ای از یک سبد خرید مشتری یا به عبارت دیگر یک تراکنش نشان داده شده است:

جدول ۱-۲ نمونه‌ای از یک تراکنش [۲]

TID (شناسه‌ی یکتا)	تراکنش (سبد خرید مربوط به هر مشتری)
۱۰۰۰	شیر، نان، تخم مرغ

فرض کنید $I = \{i_1, i_2, \dots, i_m\}$ ، مجموعه‌ی اقلام داده و m ، تعداد اقلام داده باشد. D مجموعه‌ی داده‌های تراکنشی می‌باشد که در آن برای هر تراکنش T ، یک شناسه‌ی یکتا TID وجود دارد و $T \subseteq I$ است. تراکنش T شامل مجموعه‌ی اقلام A می‌باشد اگر و فقط اگر $A \subseteq T$ باشد. یک قانون انجمنی به شکل $A \cap B = \emptyset$ و $B \subset I$ $A \subset I$ است. با داشتن مجموعه‌ی تراکنش‌ها، مسئله‌ی کشف قوانین انجمنی شامل تولید قوانین انجمنی می‌باشد که پشتیبانی و اطمینان بالاتر از کمینه‌ی پشتیبان و اطمینان تعیین شده را داشته باشند.

معیارهای مختلفی برای ارزیابی کیفیت یک قانون وجود دارد که بر اساس آن قوانین قوی، از میان مجموعه‌ی وسیعی از قوانین ممکن انتخاب می‌شوند. از معروف‌ترین و پرکاربردترین این معیارها، پشتیبانی^۱ و اطمینان^۲ می‌باشد که به صورت درصد بیان می‌شود. پشتیبانی قانون $A \rightarrow B$ عبارت است از نسبت تعداد تراکنش‌هایی که شامل A و B هستند به تعداد کل تراکنش‌ها. به عبارت دیگر احتمال آن که A و B با هم در کل تراکنش‌ها رخ دهند. پشتیبانی، فراوانی کل قانون را با توجه به کل تراکنش‌ها بیان می‌کند:

$$\text{Support}(A \rightarrow B) = P(A \cup B) = \frac{A \cup B}{D} \quad (1-2)$$

اطمینان قانون $A \rightarrow B$ عبارت است از نسبت تعداد تراکنش‌هایی که شامل A و B هستند به تراکنش‌هایی که شامل A هستند. اطمینان، قدرت دلالت قانون را بیان می‌کند:

^۱ support
^۲ confidence

$$\text{Confidence}(A \rightarrow B) = P(A|B) = \frac{A \cup B}{A} \quad (2-2)$$

۲-۲ ورودی و خروجی‌های مسئله

داده‌های ورودی برای یک الگوریتم کشف قوانین انجمنی، مجموعه‌ای از تراکنش‌ها می‌باشد که هر تراکنش زیرمجموعه‌ای از اقلام می‌باشد و شامل یک شناسه‌ی یکتا است. با داشتن مجموعه‌ی تراکنش‌ها، خروجی مسئله‌ی کشف قوانین انجمنی شامل قوانین انجمنی می‌باشد که پشتیبانی و اطمینان بالاتر از کمینه‌ی پشتیبان و کمینه‌ی اطمینان تعیین شده را داشته باشند.

۳-۲ فرض‌های مسئله

مجموعه‌ی داده‌های تراکنشی، به شکل دودویی می‌باشد. به این ترتیب که اگر قلم موردنظر در تراکنش موجود باشد، مقدار آن قلم در تراکنش مربوطه ۱ و در غیر این صورت ۰ خواهد بود. هدف کشف قوانین انجمنی مثبت می‌باشد.

۴-۲ هدف

از مشکلات مهم روش‌های سنتی برای کشف قوانین انجمنی تولید بسیار زیاد مجموعه‌های پرتکرار و پیمایش زیاد داده‌ها می‌باشد. با توجه به حجم بالای داده‌ها در مسائل مربوط به داده‌کاوی، در بیشتر موارد تمامی داده‌ها قابل بارگذاری در حافظه‌ی اصلی نمی‌باشند. از طرفی تولید مجموعه‌های پرتکرار دارای هزینه محاسباتی نمایی می‌باشد. بنابراین ارائه‌ی روش‌های مناسب برای کاهش این هزینه‌ها، حائز اهمیت می‌باشد.

۵-۲ نوآوری

در الگوریتم پیشنهادی برای کشف قوانین انجمنی یک روش کدگذاری سریع و ساده و همچنین دو روش برای عمل گر انقلاب ارائه شده است. یکی از محدودیت‌های بیشتر الگوریتم‌های تکاملی ارائه شده برای کشف قوانین انجمنی، نحوه‌ی کد کردن قوانین است که برای هر قلم موجود در

مجموعه‌ی اقلام یک یا چند مکان در نظر گرفته می‌شود و نیز ترتیب اقلام مهم می‌باشد. در این روش کدگذاری، تنها برای اقلام موجود در قانون مکان در نظر گرفته شده و هم‌چنین ترتیب اقلام مهم نیست. هم‌چنین با استفاده از دو روش ارائه شده برای عمل‌گر انقلاب، با جستجوی سریع نقاط مختلف فضای بزرگ مسئله، تعداد بیشتری قوانین در زمان کم تولید می‌گردد. راه‌کار پیشنهادی علاوه بر قوانین انجمنی، توانایی بالایی در کشف قوانین انجمنی نادر دارد.

۶-۲ سوالات تحقیق

- آیا روش خوشبندی جدید تراکنش‌ها، می‌تواند تعداد پیمایش داده‌ها را کاهش دهد؟
- آیا الگوریتم تکاملی می‌تواند قوانین انجمنی مناسب را در زمان مناسب تولید نماید؟
- آیا استفاده از الگوریتم تکاملی به طور جداگانه و موازی در هر یک از خوشبندی‌ها می‌تواند قوانین انجمنی را در زمان مناسب تولید نماید؟

۷-۲ فرضیات تحقیق

- روش خوشبندی جدید تراکنش‌ها، تعداد پیمایش داده‌ها را کاهش خواهد داد.
- الگوریتم تکاملی، قوانین انجمنی مناسب را در زمان مناسب تولید خواهد کرد.
- استفاده از الگوریتم تکاملی به طور جداگانه و موازی در هر یک از خوشبندی‌ها، قوانین انجمنی را در زمان مناسب تولید خواهد نمود.

۸-۲ خلاصه فصل

کشف قوانین انجمنی به دلیل درک راحت‌تر، یکی از مهم‌ترین شیوه‌های داده‌کاوی می‌باشد که با استفاده از این شیوه می‌توان روابط و وابستگی‌های مفیدی که در مجموعه‌های بزرگ از اقلام داده موجود می‌باشد را کشف نمود. فرآیند کشف قوانین انجمنی در روش‌های سنتی بسیار هزینه‌بر می‌باشد از طرفی حجم داده‌ها در مسائل مربوط به داده‌کاوی بالا می‌باشد، بنابراین ارائه‌ی روش‌هایی کاهش این هزینه‌ها در فرآیند کشف قوانین انجمنی بسیار مفید و موثر خواهد بود.