



دانشگاه فردوسی مشهد  
دانشکده هنری - کروه هنری کامپیوتر

پایان نامه کارشناسی ارشد

## حفظ حریم خصوصی در کاوش سودمندی

: نگارنده

طاهر ره‌گوی

استاد راهنما:

دکتر رضا منصفی

## چکیده

یکی از مسائل مهم دنیای امروز استخراج دانش از پایگاهداده‌های بسیار بزرگ است. داده‌کاوی این امکان را فراهم کرده است که بتوان به صورت خودکار اطلاعات مفیدی را از پایگاهداده‌های بسیار بزرگ استخراج نمود. اطلاعات استخراج شده توسط داده‌کاوی ممکن است موجب نقض حریم خصوصی افراد و سازمان‌ها گردد. با افزایش موارد نقض حریم خصوصی توسط داده‌کاوی، نگرانی‌ها در میان شهروندان برای حریم خصوصی نیز افزایش روزافزون داشته است. لذا در سال‌های اخیر حفظ حریم خصوصی در داده‌کاوی به عنوان مبحثی حیاتی، مورد توجه جامعه علمی قرار گرفته است. یکی از روش‌های بسیار مهم در داده‌کاوی، کاوش قوانین انجمنی می‌باشد. مبحث حفظ حریم خصوصی در کاوش سودمندی که یکی از مدل‌های کاوش قوانین انجمنی است، اخیراً در محافل علمی مطرح گردیده و تاکنون دو الگوریتم مکاشفه‌ای برای آن ارائه شده است. اگرچه این الگوریتم‌ها به خوبی می‌توانند حریم خصوصی داده‌ها را حفظ کنند اما داده‌های تولید شده توسط این الگوریتم‌ها از کیفیت مناسبی برخوردار نیستند.

از این رو هدف اصلی این تحقیق بر روی طراحی و پیاده‌سازی روش‌ها و الگوریتم‌های جدید و کارآمد برای حفظ حریم خصوصی در کاوش سودمندی مرکز دارد. مسئله یافتن راه حل بهینه برای پاکسازی پایگاهداده از مجموعه-اقلام حساس، یک مسئله برنامه‌ریزی غیرخطی عدد صحیح است و راه حل کارآمد ریاضی برای حل آن وجود ندارد. روش‌های پیشنهادی برای حل مسئله فوق در سه دسته‌ی مکاشفه‌ای، تکاملی و دقیق ارائه شده‌اند. الگوریتم‌های مکاشفه‌ای بر اساس قوانین شهودی کلی می‌کوشند که فرآیند پاکسازی را با وارد نمودن کمترین آسیب به کیفیت داده‌ها، به اتمام برسانند. الگوریتم‌های تکاملی پیشنهادی تلاش می‌کنند که با بهینه‌سازی معیارهای کیفیت داده عملیات پاکسازی را انجام دهند. در نهایت الگوریتم‌های دقیق، مدل ساده شده‌ای از مسئله فوق که با روش‌های ریاضی قابل حل باشد را به عنوان تقریبی از مسئله اصلی در نظر گرفته و با حل آن امید دارند که عملیات پاکسازی را با تولید داده‌های با کیفیت بالا به انجام برسانند. نتایج آزمایشات ارائه شده در این تحقیق، برتری چشم‌گیر روش‌های پیشنهادی بر روش‌های موجود کنونی را نشان می‌دهد.

**کلمات کلیدی:** حفظ حریم خصوصی، کاوش سودمندی، کاوش قوانین انجمنی، داده‌کاوی، برنامه‌ریزی عدد صحیح، الگوریتم‌های مکاشفه‌ای، الگوریتم‌های تکاملی

## فهرست مطالب

ii.....	چکیده
iii.....	فهرست مطالب
vii .....	فهرست شکل‌ها
ix.....	فهرست جداول
11.....	فصل ۱- مقدمه
۱۲.....	۱-۱- داده کاوی
۱۳.....	۱-۱-۱- کاوش قوانین انجمنی
۱۳.....	۱-۲- حفظ حریم خصوصی در داده کاوی
۱۴.....	۱-۲-۱- شاخه‌های تحقیقاتی در PPDM
۱۶.....	۱-۲-۲- الگوریتم‌ها و روش‌های PPDM
۱۸.....	۱-۳- مخفی‌سازی قوانین انجمنی
۱۹.....	۱-۴- حفظ حریم خصوصی در کاوش سودمندی
۲۰.....	۱-۵- نمای کلی پایان‌نامه
۲۱.....	فصل ۲- مخفی‌سازی قوانین انجمنی
۲۲.....	۲-۱- پیش‌نیاز
۲۳.....	۲-۱-۱- مقدمات و تعاریف
۲۵.....	۲-۱-۲- تئوری مرزها
۲۶.....	۲-۱-۳- تعریف رسمی مسأله
۲۶.....	۲-۳-۱-۱- اهداف متداول‌ترین‌های مخفی‌سازی قوانین انجمنی
۲۸.....	۲-۳-۱-۲- بیان مسأله
۲۹.....	۲-۳-۳- گونه ۱: مخفی‌سازی قوانین انجمنی حساس
۲۹.....	۲-۴-۳- گونه ۲: مخفی‌سازی مجموعه‌اقلام حساس

۳۰	۲-۲- دسته‌بندی متداول‌وزیری‌های مخفی‌سازی قوانین انجمنی
۳۱	۳-۲- الگوریتم‌های مکاشفه‌ای
۳۲	۱-۳-۲- الگوریتم‌های مبتنی بر درهم‌سازی
۳۴	۲-۳-۲- الگوریتم‌های مبتنی بر مسدودسازی
۳۵	۴-۲- الگوریتم‌های مبتنی بر اصلاح مرز
۳۵	۲-۴-۲- الگوریتم <i>BBA</i>
۳۶	۳-۴-۲- الگوریتم <i>Max-Min</i>
۳۶	۲-۵- الگوریتم‌های دقیق
۳۷	۱-۵-۲- الگوریتم <i>Menon</i>
۳۹	۶-۲- خلاصه
۴۰	<b>فصل ۳- حفظ حریم خصوصی در کاوش سودمندی</b>
۴۱	۱-۳- کاوش سودمندی
۴۴	۲-۳- حفظ حریم خصوصی در کاوش سودمندی ( <i>PPUM</i> )
۴۵	۱-۲-۳- فرآیند پاک‌سازی
۴۶	۲-۲-۳- محاسبه کارایی
۴۷	۳-۲-۳- الگوریتم‌های <i>PPUM</i>
۴۷	۱-۳-۲-۳- الگوریتم مخفی‌سازی قلم با سودمندی بالا، اول ( <i>HHUIF</i> )
۴۹	۲-۳-۲-۳- الگوریتم مخفی‌سازی قلم با بیشترین تداخل با مجموعه اقلام حساس، اول ( <i>MSICF</i> )
۵۰	۳-۳- خلاصه
۵۱	<b>فصل ۴- روش‌های پیشنهادی</b>
۵۶	۱-۴- مشکلات و چالش‌های موجود در حفظ حریم خصوصی در کاوش سودمندی
۵۹	۲-۴- روش‌های پیشنهادی

۱-۲-۴- رهیافت مکاشفه‌ای .....	۶۰
۱-۲-۴-۱- الگوریتم مخفی‌سازی تراکنش با بیشترین تداخل بامجموعه-اقلام حساس، اول .....(HTMSCF)	۶۲
۱-۲-۴-۲- الگوریتم مخفی‌سازی تراکنش با کمترین تداخل با مجموعه-اقلام مجاز، اول .....(HTMLCF)	۶۴
۱-۲-۴-۳- الگوریتم مخفی‌سازی تراکنش با بیشترین تداخل بامجموعه-اقلام حساس و کمترین تداخل با مجموعه-اقلام مجاز، اول (HTMMCF)	۶۵
۱-۲-۴-۴- الگوریتم مخفی‌سازی قلم با کمترین تداخل با مجموعه-اقلام مجاز، اول .....(HIMLCF)	۶۶
۱-۲-۴-۵- الگوریتم مخفی‌سازی قلم با بیشترین تداخل بامجموعه-اقلام حساس و کمترین تداخل بامجموعه-اقلام مجاز، اول (HIMMCF)	۶۷
۱-۲-۵- جمع‌بندی الگوریتم‌های مکاشفه‌ای .....	۶۷
۲-۲-۴- رهیافت تکاملی .....	۶۹
۲-۳-۴- رهیافت دقیق .....	۷۱
۳-۴- پیچیدگی زمانی الگوریتم‌های پیشنهادی .....	۷۶
۳-۴-۱- افزایش سرعت جستجو با استفاده از فایل وارونه .....	۷۷
۳-۴-۲- پیچیدگی زمانی بخش‌های مختلف الگوریتم‌ها .....	۷۷
۳-۴-۳- پیچیدگی زمانی الگوریتم‌های مکاشفه‌ای .....	۷۸
۳-۴-۴- پیاده‌سازی روش‌های پیشنهادی .....	۸۰
۴-۴- مجموعه داده‌ها .....	۸۱
۴-۵- معیارهای ارزیابی .....	۸۳
۴-۶- ارزیابی و بررسی نتایج روش‌های پیشنهادی .....	۸۵
۴-۶-۱- بهبود تابع ارزیابی الگوریتم ژنتیک .....	۸۵
۴-۶-۲- تعیین مقادیر پارامترهای الگوریتم ژنتیک .....	۸۶

۳-۶-۴- تعیین مقادیر حداقل آستانه سودمندی مناسب جهت ارزیابی الگوریتم‌ها	۹۰
۳-۶-۴- بررسی الگوریتم‌ها از نظر $HF$	۹۲
۴-۶-۴- بررسی الگوریتم‌ها از نظر $MC$	۹۲
۴-۶-۴-۱- پایگاهداده‌های کوچک	۹۲
۴-۶-۴-۲- پایگاهداده‌های بزرگ	۹۵
۴-۶-۴-۳- بررسی الگوریتم‌ها از نظر $DBDR$	۹۷
۴-۶-۴-۴-۱- پایگاهداده‌های کوچک	۹۷
۴-۶-۴-۴-۲- پایگاهداده‌های بزرگ	۱۰۰
۴-۶-۴-۴-۳- بررسی الگوریتم‌ها از نظر $TA$	۱۰۱
۴-۶-۴-۴-۱- پایگاهداده‌های کوچک	۱۰۱
۴-۶-۴-۴-۲- پایگاهداده‌های بزرگ	۱۰۴
۴-۶-۴-۴-۳- بررسی الگوریتم‌ها از دیدگاه زمان اجرا	۱۰۵
۴-۶-۴-۴-۱-۱- پایگاهداده‌های کوچک	۱۰۵
۴-۶-۴-۴-۲-۱- پایگاهداده‌های بزرگ	۱۰۸
۴-۶-۴-۴-۳- بررسی کلی الگوریتم‌ها	۱۱۰
۴-۶-۴-۴-۴- خلاصه	۱۱۲
<b>فصل ۵- نتیجه‌گیری و توصیه‌های آتی</b>	<b>۱۱۴</b>
۱-۵- کارهای آتی	۱۱۷
۱-۱-۵- استفاده از روش مسدودسازی به جای درهم‌سازی در الگوریتم‌های پیشنهادی ..	۱۱۷
۱-۲-۵- افراز پایگاهداده به زیرمجموعه‌های مستقل	۱۱۸
<b>مراجع</b>	<b>۱۱۹</b>

## فهرست شکل‌ها

شکل ۲-۱: پایگاهداده $D$ به همراه مجموعه-اقلام وقوانین انجمنی آن ..... ۲۵
شکل ۲-۲: مثالی از دو شبکه برای پایگاهداده‌هایی با $I = \{a, b, c, d\}$ (a) و $I = \{a, b, c\}$ ..... ۲۶
شکل ۲-۳: طبقه‌بندی رهیافت‌های مخفی‌سازی قوانین انجمنی از چهار دیدگاه مختلف ..... ۳۰
شکل ۳-۱: فرآیند پاک‌سازی در $PPUM$ ..... ۴۶
شکل ۴-۱: شبکه مجموعه-اقلام پایگاهداده مثال ۱ ..... ۵۳
شکل ۴-۲: دسته‌بندی الگوریتم‌های مکافهای ..... ۶۹
شکل ۴-۳: فرآیند پاک‌سازی در الگوریتم $PPUM\_GA$ ..... ۷۱
شکل ۴-۴: درصد کاهش به صفر در پایگاهداده آزمایشی ..... ۷۲
شکل ۴-۵: فرآیند پاک‌سازی در رهیافت دقیق ..... ۷۶
شکل ۴-۶: توابع توزیع $lognormal$ با $\mu = 0$ و $\delta = 1$ ..... ۸۱
شکل ۴-۷: نمودار تغییرات $MC$ برای مقادیر مختلف نرخ جهش و تقاطع ..... ۸۷
شکل ۴-۸: نمودار تغییرات $DBDR$ برای مقادیر مختلف نرخ جهش و تقاطع ..... ۸۷
شکل ۴-۹: نمودار تغییرات ترکیب $MC$ و $DBDR$ برای مقادیر مختلف نرخ جهش و تقاطع ..... ۸۸
شکل ۴-۱۰: نمودار تغییرات $MC$ برای مقادیر مختلف نرخ جهش و تقاطع ..... ۸۸
شکل ۴-۱۱: نمودار تغییرات $DBDR$ برای مقادیر مختلف نرخ جهش و تقاطع ..... ۸۹
شکل ۴-۱۲: نمودار تغییرات ترکیب $MC$ و $DBDR$ برای مقادیر مختلف نرخ جهش و تقاطع ..... ۸۹
شکل ۴-۱۳: نمودار مقایسه الگوریتم‌ها از نظر $MC$ در پایگاهداده $db1$ ..... ۹۳
شکل ۴-۱۴: نمودار مقایسه الگوریتم‌ها از نظر $MC$ در پایگاهداده $db2$ ..... ۹۳
شکل ۴-۱۵: نمودار مقایسه الگوریتم‌ها از نظر $MC$ در پایگاهداده $db3$ ..... ۹۴
شکل ۴-۱۶: نمودار مقایسه الگوریتم‌ها از نظر $MC$ در پایگاهداده $db4$ ..... ۹۴
شکل ۴-۱۷: نمودار مقایسه الگوریتم‌ها از نظر $MC$ در پایگاهداده $db5$ ..... ۹۴
شکل ۴-۱۸: نمودار مقایسه الگوریتم‌ها از نظر $MC$ در پایگاهداده $db6$ ..... ۹۶
شکل ۴-۱۹: نمودار مقایسه الگوریتم‌ها از نظر $MC$ در پایگاهداده $db7$ ..... ۹۶
شکل ۴-۲۰: نمودار مقایسه الگوریتم‌ها از نظر $DBDR$ در پایگاهداده $db1$ ..... ۹۷
شکل ۴-۲۱: نمودار مقایسه الگوریتم‌ها از نظر $DBDR$ در پایگاهداده $db2$ ..... ۹۷
شکل ۴-۲۲: نمودار مقایسه الگوریتم‌ها از نظر $DBDR$ در پایگاهداده $db3$ ..... ۹۸

شکل ۴-۲۳: نمودار مقایسه الگوریتم‌ها از نظر $DBDR$ در پایگاهداده $db4$	۹۸
شکل ۴-۲۴: نمودار مقایسه الگوریتم‌ها از نظر $DBDR$ در پایگاهداده $db5$	۹۸
شکل ۴-۲۵: نمودار مقایسه الگوریتم‌ها از نظر $DBDR$ در پایگاهداده $db6$	۱۰۰
شکل ۴-۲۶: نمودار مقایسه الگوریتم‌ها از نظر $DBDR$ در پایگاهداده $db7$	۱۰۰
شکل ۴-۲۷: نمودار مقایسه الگوریتم‌ها از نظر $TA$ در پایگاهداده $db1$	۱۰۲
شکل ۴-۲۸: نمودار مقایسه الگوریتم‌ها از نظر $TA$ در پایگاهداده $db2$	۱۰۲
شکل ۴-۲۹: نمودار مقایسه الگوریتم‌ها از نظر $TA$ در پایگاهداده $db3$	۱۰۲
شکل ۴-۳۰: نمودار مقایسه الگوریتم‌ها از نظر $TA$ در پایگاهداده $db4$	۱۰۲
شکل ۴-۳۱: نمودار مقایسه الگوریتم‌ها از نظر $TA$ در پایگاهداده $db5$	۱۰۳
شکل ۴-۳۲: نمودار مقایسه الگوریتم‌ها از نظر $TA$ در پایگاهداده $db6$	۱۰۴
شکل ۴-۳۳: نمودار مقایسه الگوریتم‌ها از نظر $TA$ در پایگاهداده $db7$	۱۰۴
شکل ۴-۳۴: نمودار مقایسه الگوریتم‌ها از نظر زمان اجرا در پایگاهداده $db1$	۱۰۶
شکل ۴-۳۵: نمودار مقایسه الگوریتم‌ها از نظر زمان اجرا در پایگاهداده $db2$	۱۰۶
شکل ۴-۳۶: نمودار مقایسه الگوریتم‌ها از نظر زمان اجرا در پایگاهداده $db3$	۱۰۶
شکل ۴-۳۷: نمودار مقایسه الگوریتم‌ها از نظر زمان اجرا در پایگاهداده $db4$	۱۰۶
شکل ۴-۳۸: نمودار مقایسه الگوریتم‌ها از نظر زمان اجرا در پایگاهداده $db5$	۱۰۷
شکل ۴-۳۹: نمودار مقایسه الگوریتم‌ها از نظر زمان اجرا در پایگاهداده $db6$	۱۰۹
شکل ۴-۴۰: نمودار مقایسه الگوریتم‌ها از نظر زمان اجرا در پایگاهداده $db7$	۱۰۹

## فهرست جداول

جدول ۳-۱: (a) یک مثال از پایگاهداده (b) جدول ارزش خارجی اقلام	۴۴
جدول ۳-۲: الگوریتم مخفی‌سازی قلم با سودمندی بالا، اول	۴۸
جدول ۳-۳: الگوریتم مخفی‌سازی قلم با بیشترین تداخل با مجموعه اقلام حساس، اول	۴۹
جدول ۴-۱: جدول (الف) تراکنش‌ها و (ب) ارزش خارجی اقلام یک پایگاهداده فرضی	۵۳
جدول ۴-۲: جداول تراکنش‌های پایگاهداده‌های پاکسازی شده	۵۴
جدول ۴-۳: کارایی جواب‌های ارائه شده در مثال ۱.	۵۵
جدول ۴-۴: شبه کد الگوریتم <i>HTMSCF</i>	۶۳
جدول ۴-۵: شبه کد الگوریتم <i>HTMLCF</i>	۶۴
جدول ۴-۶: شبه کد الگوریتم <i>HTMMCF</i>	۶۵
جدول ۴-۷: شبه کد الگوریتم <i>HIMLCF</i>	۶۶
جدول ۴-۸: شبه کد الگوریتم <i>HIMMCF</i>	۶۸
جدول ۴-۹: مشخصات پایگاهداده‌های آزمایشی	۸۳
جدول ۴-۱۰: تعداد مجموعه اقلام چندتایی در پایگاهداده‌های آزمایشی	۹۱
جدول ۴-۱۱: جدول مقایسه الگوریتم‌ها از نظر $MC(\%)$ در پایگاهداده‌های $db1$ , $db2$ و $db3$	۹۴
جدول ۴-۱۲: جدول مقایسه الگوریتم‌ها از نظر $MC(\%)$ در پایگاهداده‌های $db4$ و $db5$	۹۵
جدول ۴-۱۳: جدول تعداد رتبه‌های نخست تا سوم کسب شده توسط الگوریتم‌ها در پایگاهداده‌های کوچک از نظر $MC$	۹۵
جدول ۴-۱۴: جدول مقایسه الگوریتم‌ها از نظر $MC(\%)$ در پایگاهداده‌های $db6$ و $db7$	۹۶
جدول ۴-۱۵: جدول تعداد رتبه‌های نخست تا سوم کسب شده توسط الگوریتم‌ها در پایگاهداده‌های بزرگ از نظر $MC(\%)$	۹۷
جدول ۴-۱۶: جدول مقایسه الگوریتم‌ها از نظر <i>DBDR</i> در پایگاهداده‌های $db1$ , $db2$ و $db3$	۹۹
جدول ۴-۱۷: جدول مقایسه الگوریتم‌ها از نظر <i>DBDR</i> در پایگاهداده‌های $db4$ و $db5$	۹۹
جدول ۴-۱۸: جدول تعداد رتبه‌های نخست تا سوم کسب شده توسط الگوریتم‌ها در پایگاهداده‌های کوچک از نظر <i>DBDR</i>	۹۹
جدول ۴-۱۹: جدول مقایسه الگوریتم‌ها از نظر <i>DBDR</i> در پایگاهداده‌های $db6$ و $db7$	۱۰۰

جدول ۴-۲۰: جدول تعداد رتبه‌های نخست تا سوم کسب شده توسط الگوریتم‌ها در پایگاهداده‌های بزرگ از نظر DBDR	۱۰۱
جدول ۴-۲۱: جدول مقایسه‌ی الگوریتم‌ها از نظر $TA(\%)$ در پایگاهداده‌های $db1, db2, db3$ و $db4$	۱۰۳
جدول ۴-۲۲: جدول مقایسه‌ی الگوریتم‌ها از نظر $TA(\%)$ در پایگاهداده‌های $db4, db5$ و $db6$	۱۰۳
جدول ۴-۲۳: جدول تعداد رتبه‌های نخست تا سوم کسب شده توسط الگوریتم‌ها در پایگاهداده‌های کوچک از نظر $TA$	۱۰۴
جدول ۴-۲۴: جدول مقایسه‌ی الگوریتم‌ها از نظر $TA(\%)$ در پایگاهداده‌های $db6$ و $db7$	۱۰۵
جدول ۴-۲۵: جدول تعداد رتبه‌های نخست تا سوم کسب شده توسط الگوریتم‌ها در پایگاهداده‌های بزرگ از نظر $TA$	۱۰۵
جدول ۴-۲۶: جدول مقایسه‌ی الگوریتم‌ها از نظر زمان در پایگاهداده‌های $db1, db2, db3$ و $db4$	۱۰۷
جدول ۴-۲۷: جدول مقایسه‌ی الگوریتم‌ها از نظر زمان در پایگاهداده‌های $db4, db5$ و $db6$	۱۰۸
جدول ۴-۲۸: جدول تعداد رتبه‌های نخست تا سوم کسب شده توسط الگوریتم‌ها در پایگاهداده‌های کوچک از نظر زمان	۱۰۸
جدول ۴-۲۹: جدول مقایسه‌ی الگوریتم‌ها از نظر زمان در پایگاهداده‌های $db6$ و $db7$	۱۰۹
جدول ۴-۳۰: جدول تعداد رتبه‌های نخست تا سوم کسب شده توسط الگوریتم‌ها در پایگاهداده‌های بزرگ از نظر زمان	۱۱۰
جدول ۴-۳۱: جدول مقایسه‌ی کلی الگوریتم‌ها	۱۱۲

# فصل ۱

## مقدمه

پیشرفت‌های شایانی که در دهه‌های پیشین در زمینه سخت‌افزار و ذخیره و بازیابی اطلاعات صورت گرفته، بستری فراهم آورده است که بتوان حجم عظیمی از اطلاعات را با هزینه‌ای بسیار کم ذخیره و نگهداری کرد. بر همین اساس سازمان‌ها و شرکت‌ها توانسته‌اند حجم بسیار زیادی از تراکنش‌های مربوط به اطلاعات شخصی افراد و شرکت‌های دیگر را در انبارداده‌های<sup>۱</sup> خود ذخیره‌سازی و نگهداری کنند تا بتوانند از این اطلاعات در جهت تصمیم‌گیری، برنامه‌ریزی، ارائه خدمات بهتر و کسب سود بیشتر استفاده نمایند.

## ۱-۱- داده‌کاوی<sup>۲</sup>

حجم انبوه داده‌های جمع‌آوری شده، چالش جدیدی برای پردازش و استفاده مؤثر از داده‌ها پیش روی سازمان‌ها و مؤسسات می‌گذارد که از این چالش با نام «انفجار اطلاعات<sup>۳</sup>» یاد می‌شود. الگوریتم‌ها و روش‌های بسیار زیادی توسط محققین برای حل این معضل ارائه گردیده و بسیاری از آن‌ها با موفقیت در صنعت و تجارت به کارگرفته شده‌اند. الگوریتم‌ها و روش‌های فوق در محافل علمی به عنوان شاخه‌ای جدید به نام «داده‌کاوی» شناخته می‌شوند.

روش‌های داده‌کاوی همانند طبقه‌بندی<sup>۴</sup>، خوش‌بندی<sup>۵</sup>، کاوش قوانین انجمنی<sup>۶</sup> و ... توانسته‌اند به خوبی الگوها و نظام‌های موجود در داده‌های بسیار عظیم را استخراج کرده و آن‌ها را به صورت مدل‌ها و قوانین کلی بیان نمایند و به این ترتیب حجم عظیم داده‌ها را به مجموعه‌ای از مدل‌ها و قوانین کلی خلاصه کنند. قوانین و مدل‌های استخراج شده در اختیار مدیران و تصمیم‌گیرندگان قرار می‌گیرد تا بتوانند بر اساس آن‌ها برنامه‌ریزی و تصمیم‌سازی بهتری داشته باشند. روش‌های داده‌کاوی به مؤسسات و شرکت‌های تجاری کمک کرده است تا بتوانند خدمات مؤثرتری به مشتریان خود ارائه

<sup>1</sup>Data Warehouse

<sup>2</sup>Data Mining

<sup>3</sup>Information Overload

<sup>4</sup>Classification

<sup>5</sup>Clustering

<sup>6</sup>Association Rule Mining

نموده و در نتیجه سوددهی بیشتری داشته باشد. استفاده از روش‌های داده‌کاوی تنها به تجارت محدود نمی‌شود؛ برای مثال هر دو کاندیدای اصلی در انتخابات ریاست جمهوری آمریکا در سال ۲۰۰۴ از داده‌کاوی جهت ارزیابی و پیش‌بینی آرای خود استفاده نمودند [۳۰]. همچنین داده‌کاوی به صورت وسیع در بخش بهداشت و سلامت [۶۷]، بانکداری [۶۸]، بیمه [۶۵]، مخابرات و ارتباطات [۴۳]، کنترل کیفیت [۱۷] و ... نیز به کار برده شده است.

### ۱-۱-۱- کاوش قوانین انجمنی

یکی از روش‌های بسیار مهم در داده‌کاوی، کاوش قوانین انجمنی<sup>۱</sup> (ARM) [۵] است. کاوش قوانین انجمنی یک فرآیند دو مرحله‌ای است. در مرحله نخست مجموعه-اقلام متکرر<sup>۲</sup> کشف و استخراج می‌شوند و در مرحله دوم با استفاده از این مجموعه-اقلام، قوانین انجمنی تولید می‌شوند. تاکنون الگوریتم‌های بسیاری برای یافتن مجموعه-اقلام متکرر توسط محققان ارائه شده است. الگوریتم [۷,۹] Apriori معروف‌ترین الگوریتم در این زمینه است. یکی از مدل‌های استفاده شده در کاوش قوانین انجمنی، مدل کاوش سودمندی<sup>۳</sup> [۹] است. کاوش سودمندی کمبودهای مدل سنتی که سوددهی و تعداد کالاهای در تراکنش‌ها در نظر نمی‌گرفت را برطرف می‌سازد و قوانین مفیدتری را تولید می‌کند.

### ۱-۲- حفظ حریم خصوصی در داده‌کاوی<sup>۴</sup>

بسیاری از شرکت‌ها برای دستیابی به سود بیشتر و همچنین ارائه خدمات بهتر به کاربران خود، بخشی از اطلاعات مربوط به کاربران را با شرکت‌های همتای خود به اشتراک می‌گذارند. برخی مؤسسات دولتی نیز اطلاعات استخراج شده همانند آمارگیری‌ها و اطلاعات مربوط به بهداشت و ... را

<sup>1</sup>Association Rule Mining, ARM

<sup>2</sup>Frequent Itemsets

<sup>3</sup>Utility Mining

<sup>4</sup>Privacy Preserving Data Mining, PPDM

به صورت عمومی منتشر می‌کنند تا متخصصین و محققین بتوانند این اطلاعات را تحلیل و ارزیابی نمایند و سازمان‌ها و مؤسسه‌های دیگر نیز بتوانند از این اطلاعات برای تصمیم‌گیری و برنامه‌ریزی بهتر استفاده کنند. انتشار این نوع داده‌ها به صورت عمومی اگرچه اطلاعات خصوصی افراد و مؤسسه‌های را به صورت مستقیم افشا نمی‌کند، اما این امکان را فراهم می‌سازد که حمله‌کننده با تحلیل دقیق این داده‌ها و ترکیب آن‌ها با داده‌های عمومی دیگر، بتواند به الگوهای استنتاج‌هایی دستیابی پیدا کند که حریم خصوصی<sup>۱</sup> افراد و مؤسسه‌های را به خطر اندازد.

افزایش نگرانی‌ها در بین شهروندان و بسیاری از مؤسسه‌های برای به خطر افتادن حریم خصوصی آن‌ها موجب گردید که دولتها ادامه استفاده از روش‌های داده‌کاوی را از لحاظ قانونی مورد بازبینی قرار دهند و در مواردی استفاده از روش‌های داده‌کاوی را ناقض قوانین حفظ حریم خصوصی دانسته و آن‌ها را غیرقانونی معرفی کنند. افزایش نگرانی‌ها از یک طرف و نیاز مبرم مؤسسه‌های دولت‌ها برای استفاده از روش‌های داده‌کاوی از طرف دیگر موجب گردید که شاخه جدیدی از تحقیقات با نام حفظ حریم خصوصی در داده‌کاوی (*PPDM*) به عنوان شاخه‌ای مهم و حیاتی برای ادامه امکان استفاده از روش‌های داده‌کاوی، مورد توجه محافل علمی قرار گیرد. این مبحث نخستین بار توسط *Agrawal* و *Srikant* در سال ۲۰۰۰ به محافل علمی معرفی شد<sup>[۶]</sup> و تا کنون روش‌های زیادی برای آن ابداع گردیده و مورد استفاده قرار گرفته است.

### ۱-۲-۱- شاخه‌های تحقیقاتی در *PPDM*

شاخه‌های اصلی تحقیقات در حفظ حریم خصوصی در داده‌کاوی عبارتند از [۱]:

---

<sup>1</sup>Privacy

• **حفظ حریم خصوصی در انتشار داده‌ها<sup>۱</sup>:** در این شاخه روش‌های مختلف تبدیل، برای

تغییر داده‌ها جهت حفظ حریم خصوصی مورد بررسی قرار گرفته است. این روش‌ها شامل

روش‌های تصادفی سازی<sup>۲</sup>[۶]، گمنامی درجه- $k$ <sup>۳</sup>[۱۱,۱۲] و تنوع درجه- $l$ <sup>۴</sup>[۴۷] می‌باشد.

• **تغییر نتایج الگوریتم‌های داده‌کاوی برای حفظ حریم خصوصی:** در بسیاری از موارد

نتیجه حاصل از الگوریتم‌های داده‌کاوی مانند کاوش قوانین انجمنی و طبقه‌بندی موجب

به خطر افتادن حریم خصوصی می‌شود و نه خود داده‌ها. این امر شاخه‌ای جدید را در

PPDM به خود اختصاص داده است که هدف آن تغییر الگوریتم‌های داده‌کاوی است به

نحوی که نتایج آن‌ها حریم خصوصی را نقض نکند. به عنوان مثال در مخفی‌سازی قوانین

انجمنی، قوانینی که حریم خصوصی را به خطر می‌اندازند حذف می‌گردد.

• **بازبینی پرسش‌ها<sup>۵</sup>:** این شاخه از روش‌ها مشابه روش‌های شاخه قبلی هستند با این

تفاوت که در اینجا نتایج پرسش‌ها محدود یا دست‌کاری می‌شوند. روش‌های آشفته-

سازی نتایج پرسش‌ها در[۱۴] و روش‌های محدودسازی پرسش‌ها در[۱۵,۱۶][بررسی

شده‌اند.

• **روش‌های مبتنی بر رمزنگاری برای حفظ حریم خصوصی در سیستم‌های توزیع-**

شده: در بسیاری از حالات ممکن است که اطلاعات به صورت توزیع‌شده بر روی چندین

سایت قرار گرفته باشد. مالکان این داده‌ها ممکن است بخواهند که یک تابع مشترک<sup>۶</sup> را بر

روی همه داده‌ها اعمال کنند، به‌طوری‌که اطلاعات هر سایت برای دیگری افشا نشود و

تنها حاصل تابع برای همه آن‌ها قابل مشاهده باشد. برای انجام این کار روش‌های مختلف

<sup>1</sup>Privacy Preserving Data Publishing, PPDP

<sup>2</sup>Randomization

<sup>3</sup> $k$ -Anonymity

<sup>4</sup> $l$ -Diversity

<sup>5</sup>Query Auditing

<sup>6</sup>Perturbing

<sup>7</sup>Common Function

رمزنگاری به کار رفته است، نمونه‌هایی از کارهای انجام شده در این زمینه را می‌توان

در [۶۰] یافت.

#### • چالش‌های نظری در مسائل با ابعاد بالا(مشکل افزایش ابعاد<sup>۱</sup>): مجموعه داده‌های

واقعی معمولاً دارای ابعاد بسیار بالایی هستند، این امر الگوریتم‌های *PPDM* را چه از

لحاظ زمان محاسباتی و چه از لحاظ کارایی با مشکلاتی مواجه کرده است. برای مثال

در [۵۰] نشان داده شده است که یافتن جواب بهینه برای مسئله گمنامی درجه- $k$ ، یک

مسئله *NP-Hard* است. علاوه بر این، روش‌ها و الگوریتم‌ها هنگامی که ابعاد داده‌ها خیلی

بزرگ باشد، کارایی مناسبی از خود نشان نمی‌دهند، زیرا در این حالات ترکیب ابعاد

مختلف با داده‌های عمومی دیگر ممکن است سبب کشف اطلاعات بیشتری شود.

مجموعه‌ای از حمله‌ها به مجموعه داده‌های با ابعاد بالا را می‌توانید در [۱۹، ۲۰] بیابید.

### ۱-۲-۲-۱- الگوریتم‌ها و روش‌های *PPDM*

در این بخش برخی از روش‌ها و الگوریتم‌هایی که در هر کدام از شاخه‌های *PPDM* استفاده گردیده-

است را معرفی می‌کنیم [۲].

#### • تصادفی‌سازی<sup>۲</sup>: روش تصادفی‌سازی به صورت سنتی برای حفظ حریم خصوصی در انتشار

داده‌هایی مانند نظرسنجی‌ها استفاده می‌شود. در این روش با استفاده از یک تابع توزیع

احتمال، به داده‌ها نویز اضافه می‌گردد و به این ترتیب احتمال دسترسی به اطلاعات حساس

کاهش می‌باید. این روش در [۶] معرفی گردیده است. کمی‌سازی<sup>۳</sup> حریم خصوصی در این

روش در [۱۳, ۲۰, ۲۲] مورد بررسی قرار گرفته است. انواع حمله‌ها به روش تصادفی-

<sup>1</sup>Curse of Dimensionality

<sup>2</sup>Randomization

<sup>3</sup>Privacy Quantification

سازی نیز در [۲۱,۲۸] مورد بحث قرار گرفته‌اند. تصادفی‌سازی جریان داده‌ها<sup>۱</sup> در [۲۱,۲۸] و درهم‌سازی ضربی<sup>۲</sup> در کارهای [۱۶,۳۵,۴۴,۵۳,۵۵,۵۸] ارائه شده است. تعویض داده در [۲۳] و دیگر روش‌های تصادفی‌سازی در [۷,۲۱,۶۶,۸۰] مورد بحث قرار گرفته‌اند.

- **گمنام‌سازی گروهی<sup>۴</sup>:** روش تصادفی‌سازی، روشی بسیار ساده است که به راحتی در زمان جمع‌آوری داده‌ها قابل پیاده‌سازی است زیرا نویز اضافه شده به هر رکورد مستقل از نویز اضافه شده به دیگر رکوردها است. یکی از ضعف‌های روش تصادفی‌سازی این است که احتمال استفاده از داده‌های خارجی و ترکیب آن‌ها با داده‌های موجود را در نظر نمی‌گیرد. در [۴] نشان داده شده‌است که استفاده از رکوردهایی که به صورت عمومی در دسترس هستند و ترکیب آن‌ها با داده‌هایی که به روش تصادفی‌سازی پاک‌سازی شده‌اند، می‌تواند در موارد با بعد بالا به شدت حریم خصوصی را به خطر بیندازد. این امر در حالتی که رکورد مورد نظر به راحتی از رکوردهای محلی قابل تمییز باشد، شدیدتر است. به همین دلیل در بسیاری از رهیافت‌ها برای حفظ حریم خصوصی تلاش شده‌است که گروههایی گمنام از رکوردها ساخته شود تا حمله کننده نتواند به رکورد مورد نظر با احتمال بالا دستیابی پیدا کند. روش گمنامی درجه- $k$ -[۷۰]، تنوع درجه- $l$ -[۴۷] و نزدیکی درجه- $t^5$ -[۳۸] در این دسته قرار می‌گیرند.

- **حفظ حریم خصوصی توزیع‌شده:** هدف اصلی در بسیاری از الگوریتم‌های حفظ حریم خصوصی توزیع‌شده محاسبه توابع تجمعی آماری بر روی همه داده‌ها است به نحوی که داده‌های یک شریک برای شرکای دیگر افشا نشود. در این حالت شرکا علاقه‌مند هستند که در محاسبه توابع تجمعی آماری بر روی همه داده‌ها همکاری کنند اما به‌طور کامل به هم‌دیگر اعتماد ندارند. داده‌های به اشتراک گذاشته‌شده می‌تواند دو حالت داشته باشد: (۱)

<sup>1</sup>Data Streams

<sup>2</sup>Multiplicative Perturbation

<sup>3</sup>Data Swapping

<sup>4</sup>Group Based Anonymization

<sup>5</sup> $t$ -Closeness

داده‌ها به صورت افقی<sup>۱</sup>) به صورت عمودی<sup>۲</sup> افزار شده باشند. در افزار افقی، رکوردها بین شرکای مختلف پخش شده‌است؛ همه شرکا دارای رکوردهایی با خصیصه‌های مشابه هستند اما هر یک از شرکا بخشی از رکوردها را دارد. در افزار عمودی، خصیصه‌ها در بین شرکا پخش شده است و هریک از شرکا بخشی از خصیصه‌های هر کدام از رکوردها را در اختیار دارند. تعداد رکوردها برای همه شرکا یکسان است. اکثر روش‌های استفاده شده در این زمینه مبتنی بر رمزنگاری می‌باشند. برای اطلاعات بیشتر به [۶۰] مراجعه کنید.

#### • حفظ حریم خصوصی در نتایج الگوریتم‌های داده‌کاوی: در بسیاری از موارد، خروجی

یک الگوریتم می‌تواند توسط دشمن برای استنتاج رفتار و خصوصیات داده‌های اصلی مورد استفاده قرار گیرد. به همین دلیل بسیاری از تلاش‌ها و تحقیقات صرف محدودسازی الگوریتم‌های داده‌کاوی به منظور حفظ حریم خصوصی شده‌است. از روش‌هایی که در این دسته قرار می‌گیرند می‌توان از مخفی‌سازی قوانین انجمنی<sup>۳</sup>[۱۱]، کاهش کارایی و دقت طبقه‌بندها<sup>۴</sup>[۱۵,۵۱] و بازبینی پرسش‌ها و کنترل استنتاج‌ها<sup>۵</sup>[۱۴,۱۹] نام برد.

### ۱-۳-۱- مخفی‌سازی قوانین انجمنی

عبارةت «مخفی‌سازی قوانین انجمنی» برای نخستین بار در سال ۱۹۹۹ در [۱۱] مطرح گردید. نویسندگان این مقاله مسئله تغییر یک پایگاهداده به نحوی که برخی از قوانین انجمنی (قوانين حساس) قابل استخراج نباشند را مورد بررسی قرار دادند. آن‌ها تلاش کردند که تغییرات به نحوی باشد که کم‌ترین تأثیر را بر روی قوانین مجاز داشته باشد. آن‌ها راه حل‌های مختلفی از قبیل فازی‌سازی پایگاهداده، محدود نمودن دسترسی به داده‌ها و همچنین انتشار نمونه‌هایی از داده‌ها به جای انتشار

<sup>1</sup>Horizontally Partitioned Data

<sup>2</sup>Vertically Partitioned Data

<sup>3</sup>Association Rule Hiding

<sup>4</sup>Downgrading Classifier Effectiveness

<sup>5</sup>Query Auditing and Inference Control

کل داده‌ها را پیشنهاد کردند. به دلیل پیچیدگی زیاد این مسئله، راه حل‌های ارائه شده در آن مقاله و بسیاری از راه حل‌هایی که بعداً ارائه شدند به صورت مکاشفه‌ای بودند.<sup>۱</sup> راه حل‌های مکاشفه‌ای اگرچه به لحاظ سرعت و پیچیدگی محاسباتی<sup>۲</sup> و مقیاس‌پذیری<sup>۳</sup> از کارایی بسیار خوبی برخوردار هستند اما از مشکل نقطه بهینه محلی<sup>۴</sup> رنج می‌برند. بر همین اساس اخیراً الگوریتم‌های دقیقی برای مخفی‌سازی قوانین انجمنی معرفی شده است. این الگوریتم‌ها به ۳ دسته اساسی تقسیم‌بندی می‌شوند [۲۵]: ۱) الگوریتم‌های مکاشفه‌ای، ۲) الگوریتم‌های مبتنی بر مرز<sup>۵</sup> ۳) الگوریتم‌های دقیق.<sup>۶</sup>

#### ۴-۱- حفظ حریم خصوصی در کاوش سودمندی<sup>۷</sup>

اگرچه در سال‌های اخیر مبحث مخفی‌سازی قوانین انجمنی برای مدل سنتی (مدل کاوش اقلام<sup>۸</sup>) کاوش قوانین انجمنی<sup>[۵]</sup> به خوبی مورد بررسی قرار گرفته است اما حفظ حریم خصوصی در مدل‌های جدیدتر کاوش قوانین انجمنی مانند مدل کاوش مجموعه-اقلام<sup>۹</sup> سهم-متکرر<sup>[۱۰]</sup> و مدل کاوش سودمندی<sup>[۱۱]</sup> به دلیل پیچیدگی بیشتر نسبت به مدل سنتی، کمتر مورد توجه قرار گرفته‌اند. در سال ۲۰۱۰ مبحث حفظ حریم خصوصی در کاوش سودمندی (PPUM) توسط *Yeh* و *Hsu* مطرح گردید [۷۹]. آنها مسئله مطرح در PPUM را به صورت رسمی بیان نموده و برای آن دو الگوریتم مکاشفه‌ای جدید معرفی کردند.

<sup>1</sup>Heuristic

<sup>2</sup>Computational Complexity

<sup>3</sup>Scalability

<sup>4</sup>Local Optima

<sup>5</sup>Border Based Algorithms

<sup>6</sup>Exact Algorithms

<sup>7</sup>Privacy Preserving Utility Mining, PPUM

<sup>8</sup>Frequent Itemset Mining Model

<sup>9</sup>Share Frequent Itemset Mining Model

<sup>10</sup>Utility Mining Model

## ۱-۵- نمای کلی پایان نامه

در فصل اول مقدمه‌ای بر حفظ حریم خصوصی در کاوش سودمندی ارائه شد و دسته‌بندی‌ها و روش‌های مختلف آن به‌طور مختصر معرفی گردید. تحلیل ادبیات در قالب فصل‌های دوم و سوم ارائه شده است. در فصل دوم روش‌ها و الگوریتم‌های ارائه شده برای مخفی‌سازی قوانین انجمنی مورد بررسی قرار گرفته‌اند و دسته‌بندی‌های مختلف این روش‌ها ارائه گردیده است. در فصل سوم مسأله حفظ حریم خصوصی در کاوش سودمندی به صورت رسمی بیان شده‌است و دو الگوریتم موجود برای حل این مسأله بررسی شده است. در فصل چهارم به‌طور مفصل روش‌های پیشنهادی و همچنین نتایج ارزیابی هر کدام از آن‌ها ارائه شده است. در فصل پنجم نیز نتیجه‌گیری و توصیه‌های آتی بیان شده است.