



دانشکده مهندسی برق و کامپیوتر

پایان نامه‌ی کارشناسی ارشد در رشته‌ی مهندسی کامپیوتر - نرم افزار

طراحی و پیاده‌سازی یک الگوریتم کارآمد برای کاوش اقلام  
تکراری در جریان داده‌ها

به وسیله‌ی  
فاطمه نوری

اساتید راهنما  
دکتر محمد هادی صدرالدینی  
دکتر کورش زیارتی

شهریور ماه ۱۳۹۰



بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِيْمِ

به نام خدا

## اظهارنامه

اینجانب فاطمه نوری (۸۷۰۹۳۱) دانشجوی رشته‌ی مهندسی کامپیوتر گرایش نرم-افزار دانشکده‌ی مهندسی اظهار می‌کنم که این پایان نامه حاصل پژوهش خودم بوده و در جاهایی که از منابع دیگران استفاده کرده‌ام، نشانی دقیق و مشخصات کامل آن را نوشت‌ام. همچنان اظهار می‌کنم که تحقیق و موضوع پایان نامه‌ام تکراری نیست و تعهد می‌نمایم که بدون مجوز دانشگاه دستاوردهای آن را منتشر ننموده و یا در اختیار غیر قرار ندهم. کلیه حقوق این اثر مطابق با آیین‌نامه مالکیت فکری و معنوی متعلق به دانشگاه شیراز است.

نام و نام خانوادگی: فاطمه نوری

تاریخ و امضا: ۱۳۹۰/۰۶/۳۰

به نام خدا

طراحی و پیاده‌سازی یک الگوریتم کارآمد برای کاوش اقلام تکراری در جریان  
داده‌ها

به کوشش  
فاطمه نوری

پایان‌نامه

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

در رشته‌ی  
مهندسی کامپیوتر (نرم‌افزار)

از دانشگاه شیراز  
شیراز  
جمهوری اسلامی ایران

ارزیابی شده توسط کمیته پایان نامه با درجه: عالی  
دکتر محمدهادی صدرالدینی، دانشیار بخش مهندسی کامپیوتر (رئیس کمیته).....  
دکتر کورش زیارتی، دانشیار بخش مهندسی کامپیوتر (رئیس کمیته).....  
دکتر غلامحسین دستغیبی فرد، استادیار بخش مهندسی کامپیوتر.....  
دکتر ستار هاشمی، استادیار بخش مهندسی کامپیوتر .....  
شهریور ماه ۹۰

## تّقدیم به

دو گوهر فروزان زندگیم پدر و مادر مهربانم  
و همسر عزیزم

## سپاسگزاری

پس از سپاس به درگاه خداوند سبحان که توفيق آموختن علم را به این بنده عنایت فرمود تشکر از اسانید فرهیخته به عهده اینجانب فرض است. نخست اسانید دانشمند و گرانقدر جناب آقای دکتر محمدهادی صدرالدینی و دکتر کورش زیارتی که زحمات فراوان و بی دریغ به پای من متholm شدند و مرا با صبر و حوصله فراوانشان راهنمایی فرمودند و در طول دوران دانشجویی از فضل و دانش ایشان بهره های فراوان بردم. لذا ضمن آرزوی سلامتی و طول عمر برایشان، تشکر و امتنان قلبی خود را نسبت به این دو استاد گرانقدر ابراز می دارم و امیدوارم من نیز توانسته باشم انتظاراتشان را برآورده کنم. سپس به اسانید ارجمند و فرزانه، جناب آقای دکتر غلامحسین دستغیب و جناب آقای دکتر ستار هاشمی به خاطر اینکه تقبل نمودند که به عنوان اسانید مشاور در طول تدوین پایان نامه زحمات مرا پذیرا باشند و زحمت مطالعه پایان نامه مرا به عهده گرفتند، نهایت تشکر و قدردانی خود را اعلام می دارم و برایشان آرزوی سلامتی و طول عمر می کنم.

و همچنین از شرکت مخابرات استان شیراز به خاطر حمایت مالی که از پایان نامه اینجانب به عمل آورده، نهایت تشکر و قدردانی خود را اعلام می دارم و سرانجام، نهایت والاترین سپاس خود را از صمیم قلب به محضر خانواده عزیز و گرانقدر خود که همواره مشوق و یاری دهنده ام بودند ابراز می دارم اگر چه می دانم هیچ نوع تشکر و قدردانی، نمی تواند زحماتی را که آنان در راه به ثمر رساندن اینجانب متحمل شده اند را جبران نماید، و «اینهمه را از نظر لطف خدا می بینم».

## چکیده

### به کوشش

### فاطمه نوری

داده کاوی به معنای استخراج داده و اطلاعات غیر صریح و احتمالاً سودمندی از حجم زیادی از داده‌ها می‌باشد، که در گذشته ناشناخته و پنهان بوده‌اند. با انجام عملیات داده کاوی دانش جالب و گاه غیرمنتظره، نظم‌ها و الگوهای پنهان، یا اطلاعات سطح بالا می‌توانند از مجموعه‌ای از داده‌های مرتبط موجود در پایگاه داده استخراج شوند. یافتن اقلام تکراری یک عمل بسیار مهم در داده کاوی محسوب می‌شود. با استفاده از اقلام تکراری قوانین همبستگی استخراج و بیان می‌شوند.

یکی از بهترین الگوریتم‌های ارائه شده در این زمینه الگوریتم تک گذره به نام NewMoment می‌باشد که مجموعه‌ای از اقلام تکراری بسته را با روش پنجره‌ی کشویی حساس به تراکنش بدست می‌آورد. الگوریتم پیشنهادی ما به نام TMoment، نیز مجموعه‌ای از اقلام تکراری بسته را با روش پنجره‌ی کشویی حساس به تراکنش بدست می‌آورد، با این تفاوت که الگوریتم قبل از روش بیتی استفاده می‌کند که در پایگاه داده‌های خلوت فضای زیادی را از دست می‌دهد، اما در روش پیشنهادی، خود تراکنش‌ها نگه داشته می‌شوند، که باعث کاهش فضای مصرفی و همچنین زمان اجرا می‌شود.

در این پایان‌نامه کارایی الگوریتم جدید به همراه چند الگوریتم دیگر از لحاظ حافظه مصرفی و زمان اجرا با انجام تعدادی آزمایش مورد بررسی قرار گرفته است. ارزیابی‌های صورت گرفته نشان‌دهنده برتری الگوریتم جدید از لحاظ زمان اجرا و حافظه مصرفی است.

## فهرست مطالب

صفحة	عنوان
۱	فصل اول: کلیات
۲	۱ مقدمه
۲	۱-۱ مقدمه
۲	۱-۲ بیان مسئله
۳	۱-۳ هدف و دامنه تحقیق
۴	۱-۴ ساختار پایان نامه
۵	فصل دوم: مبانی نظری
۶	۲ مبانی نظری
۶	۱-۲ داده کاوی
۸	۱-۱-۱ مرحل داده کاوی
۱۰	۲-۲ جریان داده ای
۱۳	۱-۳-۲ مدل پردازش داده ای
۱۴	۴-۲ مدیریت حافظه
۱۵	۵-۲ ساختار داده ای فشرده
۱۶	۶-۲ کاوش الگوهای توالی
۱۷	۱-۶-۲ کاوش الگوهای توالی بسته
۱۸	۷-۲ الگوریتم های سریال برای مساله کاوش الگوهای متناوب
۱۸	۱-۷-۲ الگوریتم های سریال برای کاوش مجموعه آیتم های متناوب
۲۵	۲-۷-۲ الگوریتم های سریال برای کاوش الگوهای توالی
۳۱	فصل سوم: مروری بر کارهای پیشین
۳۲	۳ کشف اقلام تکراری در جریان داده ها
۳۴	۱-۳ الگوریتم ها براساس پنجره زمانی نشانه
۳۸	۲-۳ الگوریتم براساس پنجره کاهشی

عنوان		صفحة
۳-۳ الگوریتم‌ها براساس پنجره‌های کشویی	۳۹	
۴-۳ نخستین الگوریتم عمقی برای پیدا کردن اقلام تکراری	۴۷	
فصل چهارم: الگوریتم پیشنهادی	۵۰	
۴ الگوریتمی جدید برای کشف اقلام تکراری در جریان داده‌ها	۵۰	
۱-۴ تعریف مسئله	۵۱	
۲-۴ الگوریتم TMoment	۵۲	
۱-۲-۴ روش محاسبه پشتیبان به کمک تراکنش‌ها	۵۲	
۲-۴ ساختن TCET	۵۲	
۳-۲-۴ حذف قدیمی‌ترین تراکنش	۵۵	
۴-۲-۴ اضافه کردن تراکنش جدید	۵۶	
فصل پنجم: ارزیابی کارایی و نتیجه‌گیری	۵۹	
۵ نتیجه‌گیری و کارایی	۶۰	
۱-۵ کاوش براساس پنجره‌های مختلف	۶۱	
۲-۵ کاوش براساس اندازه پشتیبان‌های مختلف	۶۴	
۳-۵ نتیجه‌گیری	۶۶	
۳-۶ زمینه‌های تحقیقاتی آینده	۶۷	
<b>منابع</b>	۶۸	

## فهرست جدول‌ها

عنوان	صفحه
جدول ۴ - ۱) لیست تراکنش‌ها در هر پنجره ..... .....	۵۳
جدول ۵ - ۱) خصوصیات datasetها	۶۰

## فهرست شکل‌ها

صفحه	عنوان
۹	شکل ۲-۱) داده‌کاوی به عنوان یک مرحله از فرآیند کشف دانش.....
۱۸	شکل ۲-۲) نمونه‌ای از یک پایگاه توالی.....
۲۰	شکل ۲-۳) مثالی از شبکه بندی زیرمجموعه‌ای.....
۲۲	شکل ۲-۴) نمونه‌ای از یک پایگاه تراکنش.....
۲۳	شکل ۲-۵) مراحل ساخت FP-tree.....
۲۴	شکل ۲-۶) تابع FP-Growth.....
۲۵	شکل ۲-۷) مثالی از شبکه بندی زیرتوالی.....
۲۶	شکل ۲-۸) نمونه‌ای از یک پایگاه توالی.....
۲۷	شکل ۲-۹) نمونه‌ای از الگوریتم PrefixSan.....
۲۹	شکل ۲-۱۰) تابع SDB.....
۳۰	شکل ۲-۱۱) تابع BIDE.....
۳۵	شکل ۳-۱) الگوریتم IsFI-forset با آمدن abde.....
۳۵	شکل ۳-۲) الگوریتم IsFI-forset با آمدن بسته اول.....
۳۶	شکل ۳-۳) الگوریتم IsFI-forset بعد از حذف اقلام غیرتکراری.....
۳۶	شکل ۳-۴) الگوریتم IsFI-forset بعد از به پایان رسیدن پایگاه داده.....
۳۷	شکل ۳-۵) ساختار SFI-Tree برای تراکنش اول <acdf> در پنجره $W_j$ .....
۳۸	شکل ۳-۶) ساختار SFI-Tree برای پنجره W.....
۳۹	شکل ۳-۷) تغییر وزن در روش‌های مختلف.....
۴۰	شکل ۳-۸) پایگاه داده نمونه.....
۴۱	شکل ۳-۹) ساختار داده CET.....
۴۱	شکل ۳-۱۰) اضافه کردن تراکنش جدید به CET.....
۴۲	شکل ۳-۱۱) حذف تراکنش قدیمی از ساختار CET.....
۴۲	شکل ۳-۱۲) ساختار FP-Tree.....
۴۳	شکل ۳-۱۳) پایگاه داده نمونه.....
۴۳	شکل ۳-۱۴) روش نمایش بیتی برای پایگاه داده نمونه.....
۴۴	شکل ۳-۱۵) ساختار NewCET در پنجره اول.....
۴۴	شکل ۳-۱۶) حذف قدیمی ترین تراکنش در ساختار NewCET.....
۴۵	شکل ۳-۱۷) اضافه کردن تراکنش جدید در ساختار NewCET.....

شکل-۳-۱۸) نحوه ذخیرهسازی نودهای CPS-Tree	۴۶
شکل-۳-۱۹) پایگاه داده نمونه برای CPS-Tree	۴۶
شکل-۳-۲۰) ساختار CPS-Tree بعد از ورود pane اول	۴۶
شکل-۳-۲۱) ساختار CPS-Tree بعد از ورود pane دوم	۴۷
شکل-۳-۲۲) ساختار CPS-Tree در پایان پنجره اول	۴۷
شکل-۳-۲۳) پایگاه داده نمونه برای الگوریتم Eclat	۴۸
شکل-۳-۲۴) درخت اولیه الگوریتم Eclat	۴۸
شکل-۳-۲۵) درخت الگوریتم Eclat	۴۹
 شکل-۴-۱) مثالی از پنجره حساس به تراکنش	۵۱
شکل-۴-۲) ساختار TCET برای پنجره اول	۵۴
شکل-۴-۳) شبه کد TCET	۵۵
شکل-۴-۴) حذف تراکنش اول	۵۶
شکل-۴-۵) کد تابع Eliminate برای حذف تراکنش قدیمی	۵۷
شکل-۴-۶) اضافه کردن تراکنش جدید به درخت TCET	۵۷
شکل-۴-۷) کد تابع اضافه کردن تراکنش جدید	۵۸
 شکل-۵-۱) مقایسه زمان اجرا الگوریتم‌های TMoment .NewMoment .Moment	۶۲
شکل-۵-۲) مقایسه حافظه مصرفی در الگوریتم‌های TMoment .NewMoment .Moment	۶۳
شکل-۵-۳) مقایسه زمان اجرا الگوریتم‌های TMoment .NewMoment .Moment	۶۵
شکل-۵-۴) مقایسه حافظه مصرفی در الگوریتم‌های TMoment .NewMoment .Moment	۶۶

# فصل اول

## مقدمه

### ۱-۱ مقدمه

امروزه حجم زیادی از داده‌ها در پایگاه داده‌های مربوط به شرکت‌ها، مراکز تجاری و دولتی ذخیره می‌شود. استفاده معمول از این داده‌ها انجام عملیات گزارش‌گیری برای کاربران و مدیران است. استفاده دیگری که امروزه از حجم انبوه داده‌های ذخیره شده در پایگاه‌های داده و انبارهای داده می‌شود انجام عملیات داده کاوی است. در عملیات داده کاوی به دنبال الگوهای پنهان و احتمالاً سودمند هستیم. برخی از این الگوها در انجام تصمیم‌گیری‌های به مدیران کمک می‌کنند و برخی دیگری برای کاربران و مشتریان مفید واقع می‌شوند. الگوایی که در عملیات مختلف داده کاوی پیدا می‌شوند انواع گوناگونی دارند. یک نوع معروف و پرکاربرد از این الگوها قواعد وابستگی یا قوانین وابستگی نام دارد.

### ۲-۱ بیان مسئله

معروف‌ترین کاربرد قوانین وابستگی در تحلیل سبد خرید برای فروشگاه‌ها و مراکز تجاری است. به عنوان مثال پس از کاوش پایگاه داده مربوط به یک فروشگاه زنجیره‌ای ممکن است مشخص شود که مشتریانی که از این فروشگاه مربا می‌خرند به احتمال ۶۰٪ کره نیز خواهند خرید. یافتن چنین قواعدی می‌تواند در چیدن قفسه‌ها، راهنمای مشتریان و مسائل مدیریتی سودمند باشد. عملیات یافتن قواعد وابستگی را کشف یا کاوش قواعد وابستگی گویند.

داده کاوی<sup>۱</sup> به معنای استخراج داده و اطلاعات غیر صریح و احتمالاً سودمندی از حجم زیادی از داده‌ها می‌باشد، که در گذشته ناشناخته و پنهان بوده‌اند. با انجام عملیات داده کاوی دانش جالب و گاه غیرمنتظره، نظم‌ها و الگوهای پنهان، یا اطلاعات سطح بالا می‌توانند از مجموعه‌ای از داده‌های مرتبط موجود در پایگاه داده استخراج شوند و از زوایای مختلف مورد بررسی قرار گیرند. بنابراین پایگاه‌های داده حجیم را می‌توان به عنوان منابعی غنی و قابل اطمینان برای تولید و وارسی برخی اطلاعات در نظر گرفت.

<sup>۱</sup> Data mining

یافتن اقلام تکراری<sup>۱</sup> یک عمل بسیار مهم در داده کاوی محسوب می شود. با استفاده از این اقلام تکراری قوانین همبستگی استخراج و بیان می شوند. در یک دسته از تراکنش‌ها، پشتیبان<sup>۲</sup> اقلام، درصد تکرار اقلام در کل تراکنش‌هاست. یک عنصر تکراری است اگر پشتیبان آن بزرگتر یا مساوی یک مینیمم پشتیبانی شود، که کاربر در نظر گرفته است.

## ۲-۱ هدف و دامنه تحقیق

در سالهای اخیر استخراج جریان داده‌ها<sup>۳</sup> بسیار مورد توجه محققان قرار گرفته است. با به وجود آمدن محیط‌های کاربردی جدید مانند آنالیز ترافیک شبکه، کاوش جریان‌های کلیک وب<sup>۴</sup> و کشف نفوذ به شبکه جریان داده‌ها با سرعت بالا، پیوسته<sup>۵</sup> و نامحدود وارد سیستم می‌شوند و دیگر پردازش داده‌ها نمی‌تواند به صورت ایستا<sup>۶</sup> انجام شود. انواع مختلف داده‌های مربوط به سری‌های زمانی، داده‌هایی که در محیط‌های پویا تولید می‌شوند مانند میزان مصرف انرژی، ترافیک شبکه، مبادلات سهام، ارتباطات از راه دور، جریان‌های کلیک وب، مراقبت‌های بصری و امور نظارتی در مسایل محیطی و آب و هوا، مثال‌های از جریان‌های داده‌ایی محسوب می‌شوند.

در جریان داده‌ها امکان استفاده از روش‌های مشابه که از پایگاه داده‌های سنتی استفاده می‌کنند، وجود ندارد. به دلیل حجم زیاد داده‌ها در سیستم‌ها، داده‌ها را نمی‌توان به صورت ماندگار ذخیره کرد و به صورت جریان‌های موقت می‌باشند. در جریان داده‌ها هر اقلامی<sup>۷</sup> را فقط یکبار می‌توان آزمایش کرد، همچنین به خاطر نامحدود بودن داده‌ها سیستم با محدودیت حافظه روبرو است و در ضمن نتایج کاوش باید خیلی سریع تولید شود.

---

<sup>1</sup> Frequent itemset

<sup>2</sup> transaction

<sup>3</sup> support

<sup>4</sup> Data stream

<sup>5</sup> Web click stream mining

<sup>6</sup> continuous

<sup>7</sup> static

<sup>8</sup> itemsets

در جریان داده‌ها یکی از موضوع‌های مهم کاوش اقلام تکراری<sup>۱</sup> است که در سالهای اخیر بسیار مورد توجه محققان قرار گرفته است.

با توجه به اهمیت جریان داده‌ها که گفته شد در این تحقیق ما سعی برای داریم که الگوریتم جدیدی را برای کاوش اقلام تکراری در جریان داده‌ها مطرح کنیم.

### ۳-۱ ساختار پایان‌نامه

در این پایان‌نامه ابتدا به معرفی داده کاوی و جریان‌های داده‌ای می‌پردازیم. سپس از میان عملیات داده کاوی مسئله کشف قواعد داده‌ای را به طور دقیق بررسی خواهیم کرد. روش‌ها و الگوریتم‌های ارائه شده برای یافتن اقلام تکرار شونده که مهمترین گام برای کشف قواعد وابستگی است را بررسی خواهیم کرد. این الگوریتم‌ها از جنبه‌های مختلف مقایسه و مزایا و معایب هر کدام بررسی می‌شوند. در ادامه یک الگوریتم جدید برای کشف اقلام تکرار شونده در جریان داده‌ای ارائه می‌کنیم و در پایان ارزیابی و نتایج آزمایشات و نتیجه‌گیری نهایی را خواهیم داشت.

---

<sup>1</sup> Frequent itemset mining

# فصل دوم

## مبانی نظری

### ۱-۲ داده کاوی

اخیرا سرعت تولید و جمع آوری دادهها 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 دادهها ما را با اینوی از دادهها مواجه کرده است. امروزه پایگاه داده در زمینه های تجاری، اداری، علمی، مهندسی و زمینه های دیگر استفاده می شوند. تعداد چنین پایگاه دادهای به دلیل نیاز مبرم به جمع آوری و گزارش گیری از دادهها و همچنین وجود سیستم های قدرتمند پایگاه داده در حال افزایش است. این چنین رشد فزاینده ای در دادهها نیاز ضروری برای ابزار جدیدی که بتواند به طور بتواند به طور هوشمند و خودکار این دادهها پردازش کرده و به اطلاعات و دانش های سودمند تبدیل کند احساس می شود. در نتیجه داده کاوی به یک زمینه تحقیقاتی با اهمیت فراوان تبدیل شده است.

از داده کاوی همچنین به عنوان کشف دانش در بانک های دادهای یاد می شود، که به معنی فرایند استخراج اطلاعات غیرصریح و احتمالا سودمندی از پایگاه داده است که در گذشته ناشناخته و پنهان بوده اند. با انجام عملیات داده کاوی دانش های جالب و گاه غیرمنتظره، نظم ها و الگوهای پنهان یا اطلاعات سطح بالا می توانند از داده های مرتبط در پایگاه داده استخراج شوند و از زوایای مختلف مورد بررسی قرار گیرند. بنابراین پایگاه های داده حجمی را می توان به عنوان منبعی غنی و قابل اطمینان برای تولید و وارسی برخی از دانش ها و اطلاعات در نظر گرفت.

کاوش اطلاعات و دانش از پایگاه داده حجمی به عنوان یک موضوع کلیدی برای محققینی که در زمینه پایگاه های داده و یادگیری ماشین کار می کنند و به فرصتی برای کسب در آمدهای 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 تصییم 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 کاربردهای مختلف، داده کاوی در سال های اخیر توجه زیادی را به خود جلب کرده است [۱، ۲].

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، خوشه بندی، پیش بینی و

تحلیل سری‌های زمانی از جمله مهمترین روش‌ها و تکنیک‌های داده کاوی به شمار می‌آیند. در ادامه هر کدام از روش‌ها به صورت خلاصه توضیح داده می‌شود.

کشف قواعد وابستگی در جریان‌های داده‌ای از جذبیت زیادی در بین متخصصان داده کاوی برخوردار است. در این تکنیک وابستگی‌ها و ارتباطات میان داده‌ها بدست می‌آید. نتیجه این عملیات داده کاوی دسته‌ای از قواعد است که به آنها قواعد وابستگی گفته می‌شود. یکی دیگر از روش‌های مهم داده کاوی توانای انجام عملیات دسته‌بندی در حجم زیاد داده‌هاست. این عملیات کاوش قوانین دسته‌بندی نیز نامیده می‌شود. در این روش اشیا موجود در یک پایگاه داده براساس مقادیر چند خصوصیت از آنها به دسته‌های مجزا تقسیم می‌شوند. در دسته‌بندی داده‌ها آزمایشی تحلیل می‌شوند. در دسته‌بندی داده‌ها مجموعه داده‌های آزمایشی برچسب کلاس‌ها مشخص است. برای هر کلاس داده‌های آزمایشی مدلی براساس خصوصیات داده‌ها ساخته می‌شود. حاصل عملیات دسته‌بندی می‌تواند یک درخت تصمیم یا مجموعه‌ای از قوانین دسته‌بندی باشد که برای فهم بهتر داده‌های موجود در پایگاه داده و همچنین دسته‌بندی داده‌هایی که در آینده به پایگاه داده اضافه می‌شوند به کار می‌رود. دسته‌بندی داده‌ها ارتباط تنگاتنگی با کاوش قواعد وابستگی دارد، به طوری که گاهی دسته‌بندی را به کمک قواعد وابستگی انجام می‌دهند.

به عنوان مثال برای فروشنده‌ی ماشین دسته‌بندی مشتریانش براساس تمایل و علاقه هر کدام به انواع مختلف ماشین مطلوب است به طوری که بتواند به مشتریانش خدمات بهتری ارائه دهد و کاتالوگ‌های محصولات جدید مورد نظر آنها را برایش بفرستد و درآمد حاصل از فروش خود را افزایش دهد. خوشبندی داده‌ها از مهمترین روش‌های داده کاوی به شمار می‌آید. در خوشبندی مجموعه‌ای از داده‌ها گروه‌بندی می‌شوند. فرق خوشبندی با دسته‌بندی داده‌ها در این است که در خوشبندی برخلاف دسته‌بندی تعداد کلاس‌ها در ابتدا مشخص نیستند. خوشبندی داده‌ها براساس اصل مفهومی زیر صورت می‌گیرد:

حداکثر کردن شباهت‌های اعضای هر کلاس و حداقل کردن شباهت‌ها بین اعضای مربوط به کلاس‌های مختلف.

به عنوان مثالی از خوشبندی مجموعه‌ای از کالاها را می‌توان در ابتدا به صورت مجموعه‌ای از کلاس‌های مختلف خوشبندی کرده و سپس مجموعه‌ای از قوانین را براساس این چنین دسته‌بندی نتیجه‌گیری کرد.

پیش‌بینی یکی دیگر از تکنیک‌های داده کاوی است که در آن مقادیر ممکن برای متغیرهای نامعلوم پیش‌بینی می‌شوند. در پیش‌بینی ابتدا داده‌ایی که به متغیر نامعلوم مربوط هستند به وسیله برخی تحلیل‌های آماری پیدا می‌شوند، سپس از برخی روش‌های هوشمند مانند شبکه‌های عصبی و الگوریتم‌های ژنتیک برای انجام پیش‌بینی استفاده می‌شود. برای مثال مقدار حقوقی که یک کارمند می‌گیرد را می‌توان با استفاده از چگونگی توزیع حقوق کارمندان مشابه