

داده‌کاوی در آمار رسمی

مطالعه موردی: بررسی الگوی مصرف خانوارهای شهری براساس طرح هزینه و درآمد خانوار

مجری

عباس مرادی

همکاران

کاوه کیانی

اشکان شباک

ایوب فرامرزی

حامد لروند



پژوهشکده‌ی آمار

گروه پژوهشی پردازش داده‌ها و اطلاع‌رسانی

بهار ۱۳۹۸

پیش‌گفتار

عصر حاضر، عصر تغییرات سریع، پی‌درپی و متنوع است. همزمان با این تغییرات، به ویژه برخط شدن و به‌روز شدن سامانه‌ها و سیستم‌های ثبت و گردآوری داده‌ها، سازمان‌های آماری خصوصاً سازمان‌های متولی آمار رسمی را با انباشت زیاد داده‌ها و اطلاعات آماری ساختار نیافته در قالب‌ها و فرمت‌های متنوع مواجه کرده است. در واقع امروزه «انفجار داده‌ها» به وضعیتی رسیده است که با آمار کلاسیک و روش‌های سنتی تجزیه و تحلیل آن امکان‌پذیر نیست. برطرف شدن این چالش بزرگ، محققان ریاضی، آمار و علوم کامپیوتر را بر آن داشت تا راهی برای تجزیه و تحلیل این داده‌ها کشف نمایند. مهمترین یافته و راه حل این چالش برای محققان علوم داده، روش‌ها و الگوریتم‌هایی است که امروزه «داده‌کاوی» نام گرفته است. سازمان‌های متولی آمار رسمی و در راس این سازمان‌ها، مرکز آمار ایران به عنوان یگانه مرجع آمارهای رسمی کشور نیز نیازمند دانشی است که به کمک آن داده‌ها را تجزیه و تحلیل نماید و در اختیار سیاست‌گذاران و مدیران ارشد کشور به منظور تدقیق نمودن تصمیم‌گیری‌ها قرار دهد. پر واضح است که دانش و خرد برگرفته از داده‌های خام و متنوع، تصمیم‌گیری و برنامه‌ریزی بهینه برای برنامه‌های کلان و آتی کشور را تسهیل می‌نماید. امروزه مهم‌ترین و در واقع تنهاترین ابزار برای رفع این چالش، استفاده از تکنیک‌های داده‌کاوی است. با استفاده از دانش داده‌کاوی، روابط کشف و الگوهای پنهان در داده‌ها آشکار می‌شوند. در واقع داده‌کاوی به دانشی بین‌رشته‌ای (ریاضی، آمار و علوم کامپیوتر) که هدف آن ایجاد ابزارهایی در جهت پردازش مه‌داده‌ها و استخراج الگوهای مخفی و قوانین حاکم بر آنهاست، اطلاق می‌شود. داده‌کاوی، در تمامی زمینه‌هایی که داده‌ها گردآوری شوند و نیازمند تجزیه و تحلیل باشند، نقش بسزایی خواهد داشت.

پژوهشکده‌ی آمار با توجه به رسالت خود در جهت تسهیل این مهم برای تجزیه و تحلیل داده‌های آمار رسمی اجرای طرح پژوهشی «داده‌کاوی در آمار رسمی» را در دستورکار خود قرار داد. هدف این طرح پژوهشی معرفی برخی از روش‌های داده‌کاوی و کاربردهای آن در آمار رسمی است. این پژوهش در گروه پژوهشی پردازش داده‌ها و اطلاع‌رسانی پژوهشکده‌ی آمار توسط آقای دکتر عباس مرادی به عنوان مجری و آقایان دکتر ایوب فرامرزی، حامد لرونند، دکتر کاوه کیانی و دکتر اشکان شباک به‌عنوان همکار به انجام رسیده است، که از آنان صمیمانه تشکر و قدردانی می‌شود. همچنین از دکتر محمد شیرینی، به خاطر نظرات ارزشمندشان تشکر و قدردانی می‌شود.

گروه پژوهشی پردازش داده‌ها و اطلاع‌رسانی

پژوهشکده‌ی آمار

فهرست مطالب

۱۰	۱- کلیات تحقیق
۱.....	۱-۱- مقدمه
۱.....	۲-۱- بیان مسئله
۲.....	۳-۱- اهداف
۳.....	۴-۱- پیشینه تحقیق
۳.....	۵-۱- روش تحقیق
Error! Bookmark not defined.	۲- برخی از الگوریتم‌های داده‌کاوی
Error! Bookmark not defined.	۳- آمار رسمی
۲۷.....	۱-۳- مقدمه
۲۸.....	۲-۳- اصول بنیادی‌ترین آمارهای رسمی
۲۹.....	۳-۳- هزینه و درآمد خانوار
Error! Bookmark not defined.	۴- کاربردهای داده‌کای در آمار رسمی
۳۱.....	۱-۴- مقدمه
۳۲.....	۲-۴- خوشه‌بندی خانوارهای شهری بر اساس داده‌های هزینه و درآمد سال ۱۳۹۰
۴۰.....	۳-۴- وابستگی سبد هزینه‌ی خانوارهای شهری خوشه‌بندی شده سال ۱۳۹۰
۵۲.....	۴-۴- خوشه‌بندی خانوارهای شهری بر اساس داده‌های هزینه و درآمد سال ۱۳۹۶
۵۸.....	۵-۴- وابستگی سبد هزینه خانوارهای شهری خوشه‌بندی شده سال ۱۳۹۶
۶۸.....	۶-۴- نتایج
Error! Bookmark not defined.....	پیوست
۱۵۵.....	مرجع‌ها

فهرست جداول

- جدول ۱-۲ - سه مشاهده با چهار متغیر مفروض ۲۱
- جدول ۱-۴ - خصوصیات کلی خوشه‌های ایجاد شده خانوارهای شهری در سال ۱۳۹۰ ۳۷
- جدول ۲-۴ - خصوصیات کلی خوشه‌های ایجاد شده خانوارهای شهری در سال ۱۳۹۶ ۵۵
- جدول ۳-۴ - میزان انرژی دریافتی برخی از کالاها برحسب کالری ۶۸

فهرست شکل‌ها

- شکل ۱-۱- هرم دانش ۴
- شکل ۱-۲- داده‌های دانش‌آموزان شهرستان در یک نگاه ۹
- شکل ۲-۲- درخت تصمیم ۱۰
- شکل ۲-۳- تأثیرات متغیرهای مهم در ادامه تحصیل ۱۱
- شکل ۲-۴- تأثیرپذیری متغیرهای ورودی بر متغیر هدف ۱۱
- شکل ۲-۵- بیشترین تأثیرپذیری متغیرهای ورودی بر متغیر هدف ۱۲
- شکل ۲-۶- درخت تصمیم با جزئیات ۱۳
- شکل ۲-۷- نمونه‌ای از داده‌های آموزشی ۱۵
- شکل ۲-۸- درخت تصمیم در حالت کلی برای داده‌های آموزشی مثال فوق ۱۷
- شکل ۲-۹- درخت تصمیم با جزئیات ۱۸
- شکل ۲-۱۰- خوشه‌بندی نرم ۲۳
- شکل ۲-۱۱- محاسبه‌ی معیارهای قواعد پیوند ۲۵
- شکل ۴-۱- پایگاه داده رابطه‌ی هزینه و درآمد خانوار شهری سال ۱۳۹۰ ۳۳
- شکل ۴-۲- دیاگرام (تمامی) ارتباط‌های خوشه‌های بدست آمده خانوارهای شهری سال ۱۳۹۰ ۳۴
- شکل ۴-۳- دیاگرام (برخی از) ارتباط‌های خوشه‌های بدست آمده خانوارهای شهری سال ۱۳۹۰ ۳۵
- شکل ۴-۴- دیاگرام (قویترین) ارتباط‌های خوشه‌های بدست آمده خانوارهای شهری سال ۱۳۹۰ ۳۶
- شکل ۴-۵- وابستگی خرید خوشه‌ی اول سال ۱۳۹۰ در حالت کلی ۴۰
- شکل ۴-۶- وابستگی خرید خوشه‌ی اول سال ۱۳۹۰ در حالت ارتباط کمتر، با انتخاب گره دلخواه ۴۰
- شکل ۴-۷- بزرگ‌نمایی وابستگی خرید خوشه‌ی اول سال ۱۳۹۰ در حالت کلی ۴۱
- شکل ۴-۸- مشاهده بعضی از قوانین و میانگین‌گیری از بعضی از قانون‌ها ۴۲
- شکل ۴-۹- وابستگی خرید خوشه‌ی دوم سال ۱۳۹۰ در حالت کلی ۴۲
- شکل ۴-۱۰- وابستگی خرید خوشه‌ی دوم سال ۱۳۹۰ - انتخاب یک گره ۴۳
- شکل ۴-۱۱- بزرگ‌نمایی وابستگی خرید خوشه‌ی دوم سال ۱۳۹۰ ۴۳
- شکل ۴-۱۲- مشاهده بعضی از قوانین - الف ۴۴
- شکل ۴-۱۳- مشاهده بعضی از قوانین - ب ۴۴
- شکل ۴-۱۴- وابستگی خرید خوشه‌ی سوم سال ۱۳۹۰ ۴۵
- شکل ۴-۱۵- وابستگی خرید خوشه‌ی سوم سال ۱۳۹۰ انتخاب یک گره دلخواه ۴۵
- شکل ۴-۱۶- مشاهده بعضی از قوانین در خوشه‌ی سوم سال ۱۳۹۰ ۴۶

- شکل ۴-۱۷- وابستگی خرید خوشه‌ی چهارم سال ۱۳۹۰ ۴۷
- شکل ۴-۱۸- وابستگی خرید خوشه‌ی چهارم سال ۱۳۹۰ انتخاب یک گره دلخواه ۴۷
- شکل ۴-۱۹- بزرگنمایی وابستگی خرید خوشه‌ی چهارم سال ۱۳۹۰ ۴۸
- شکل ۴-۲۰- وابستگی خرید خوشه‌ی چهارم سال ۱۳۹۰ مشاهده داده‌ها ۴۸
- شکل ۴-۲۱- وابستگی خرید کنسرو ماهی و برنج ۴۹
- شکل ۴-۲۲- وابستگی خرید کره حیوانی و تخم‌مرغ ۴۹
- شکل ۴-۲۳- وابستگی خرید کره حیوانی و تخم‌مرغ (میانگین‌گیری) ۵۰
- شکل ۴-۲۴- پایگاه داده رابطه‌ی هزینه و درآمد خانوار شهری سال ۱۳۹۶ ۵۲
- شکل ۴-۲۵- خوشه‌بندی نرم خانوارهای شهری سال ۱۳۹۶ - بیشترین ارتباطات ۵۳
- شکل ۴-۲۶- خوشه‌بندی نرم خانوارهای شهری سال ۱۳۹۶ - ارتباطات کمتر الف ۵۳
- شکل ۴-۲۷- خوشه‌بندی نرم خانوارهای شهری سال ۱۳۹۶ - ارتباطات کمتر ب ۵۴
- شکل ۴-۲۸- خوشه‌بندی نرم خانوارهای شهری سال ۱۳۹۶ - کمترین ارتباطات ممکن ۵۴
- شکل ۴-۲۹- ایجاد پایگاه‌داده‌ی رابطه‌ی برای یافتن قواعد پیوند خوشه‌های چهارم الی دهم خانوارهای سال ۱۳۹۶ ۵۸
- شکل ۴-۳۰- وابستگی خرید خوشه‌ی اول سال ۱۳۹۶ ۵۸
- شکل ۴-۳۱- وابستگی خرید خوشه‌ی دوم سال ۱۳۹۶ ۵۹
- شکل ۴-۳۲- وابستگی خرید خوشه‌ی سوم سال ۱۳۹۶ ۶۰
- شکل ۴-۳۳- میانگین‌گیری بعضی از قوانین ۶۱
- شکل ۴-۳۴- وابستگی خرید خوشه‌ی چهارم سال ۱۳۹۶ ۶۲
- شکل ۴-۳۵- وابستگی خرید خوشه‌ی پنجم سال ۱۳۹۶ ۶۲
- شکل ۴-۳۶- وابستگی خرید خوشه‌ی ششم سال ۱۳۹۶ ۶۳
- شکل ۴-۳۷- وابستگی خرید خوشه‌ی هفتم سال ۱۳۹۶ ۶۴
- شکل ۴-۳۸- وابستگی خرید خوشه‌ی هشتم سال ۱۳۹۶ ۶۵
- شکل ۴-۳۹- وابستگی خرید خوشه‌ی نهم سال ۱۳۹۶ ۶۶
- شکل ۴-۴۰- وابستگی خرید خوشه‌ی دهم سال ۱۳۹۶ ۶۷



کلیات تحقیق

۱-۱- مقدمه

داده‌کاوی^۱ بین محققان علوم آمار، کامپیوتر و ریاضی در جهت تحلیل داده‌ها و پایگاه‌های داده‌ای که شامل مه‌داده‌ها^۲ و یا داده‌های با ابعاد بالا^۳ هستند، به عنوان یکی از ابزارهای پرتوان شناخته شده است. بر همین اساس مرکز آمار ایران به عنوان یگانه مرجع آمار رسمی کشور روز به روز به اهمیت داده‌کاوی پی برده و این مهم را در اولویت‌های برنامه‌ای خود قرار داده است. پرواضح است که ایجاد یک نظام کارآمد و مؤثر در تولید و عرضه آمار از الزامات اولیه و ضروری در برنامه‌ریزی است و زیربنای برنامه‌ریزی مناسب، اطلاعات جامع، منسجم و به‌روز است. داده‌کاوی که به معنای کشف دانش و استخراج از مقادیر زیادی از داده‌های خام است. از این رو داده‌کاوی مفیدترین ابزار برای استفاده از تحلیل وقایع گذشته در پایگاه‌های داده‌ای مراکز آمار ایران خواهد بود. بدین معنی برای کشف و استخراج اطلاعات از آمار رسمی کشور که آمار مد نظر، سیاست‌گذاران، برنامه‌ریزان و مدیران کشور است نیازمند استفاده از داده‌کاوی خواهیم بود. هدف این طرح پژوهشی، معرفی برخی از روش‌های سودمند داده‌کاوی در آمار رسمی خواهد بود.

۱-۲- بیان مسئله

^۱ Data Mining

^۲ Big Data

^۳ High Dimensional Data

آمار رسمی^۴ به اطلاعات عددی گفته می‌شود که توسط دولت یا مراجع صلاحیت‌دار که در قوانین و مقررات مشخص هستند، تولید و منتشر می‌شود و اطلاعاتی را در مورد وضعیت عمومی کشور برای امور مدیریتی (برنامه‌ریزی، سیاست‌گذاری، و تصمیم‌گیری) به دست می‌دهد. در واقع آمار رسمی به اطلاعات عمومی مربوط می‌شود که به نفع جامعه و با بودجه دولتی تولید می‌شود. آمار رسمی برای همگان قابل دسترس خواهد بود. از این‌رو ذی‌نفعان را قادر می‌سازد تا با اتکا بر این آمارها، تصمیم‌های مهمی در زندگی شخصی و تجاری خویش انجام دهند. آمار رسمی مطابق با رده‌بندی‌های بین‌المللی و با رویکرد اصل بی‌طرفی، قابلیت اطمینان، اثربخشی، محرمانه بودن و شفافیت تولید و منتشر می‌شود. اما از یک سو حجم این اطلاعات بسیار زیاد و همواره در حال تغییر است. بدین معنی که این داده‌ها با حجم، تنوع و سرعت بالا در حال تولید هستند. در واقع اکثر سازمان‌های آماری دارای پایگاه‌های داده‌ای هستند که داده‌های آن از نظر حجم و تنوع بسیار گوناگون است. از این‌رو می‌توان گفت این سازمان‌ها با داده‌ها و اطلاعات آماری جمع‌آوری شده که دارای رشد حیرت‌آوری است مواجه شده‌اند به گونه‌ای که تحلیل داده‌ها با روش‌های کلاسیک برای آن‌ها امکان‌پذیر نخواهد بود. از سوی دیگر با گسترش روزافزون فناوری اطلاعات و تجهیز سازمان‌ها، شرکت‌های دولتی و خصوصی امکان جمع‌آوری اطلاعات دقیق و به روز فراهم آمده است، به طوری که مجموعه وسیع و حجیمی از داده‌ها شامل داده‌های ثبتي خلق گردیده است. این اطلاعات مد نظر مدیران و پژوهشگران بوده و از آن‌ها برای تهیه گزارش‌های مختلف و تصمیم‌گیری‌های روزمره و راهبردی استفاده می‌شود. استخراج اطلاعات از این داده‌ها نیازمند کاوش و جست‌وجوی اساسی است. تصمیم‌گیری و برنامه‌ریزی بهینه توسط سیاست‌گذاران کلان یک کشور نیازمند داشتن دانش مشتق شده از اطلاعات آماری و آمار رسمی کشور است. بدیهی است که ارتقا و توسعه علمی کشور نیازمند تصمیم‌گیری بهینه از الگوهای نهان در پایگاه‌های داده‌ای کشور است. از این‌رو سازمان‌ها و شرکت‌های دولتی نیازمند به کارگیری دانش داده‌کاوی و تکنیک‌های هوشمند آن در پایگاه‌های داده‌ای خود می‌باشند. یکی از مهمترین روش‌های داده‌کاوی الگویابی و مدلسازی است. داده‌کاوی داده‌ها را به منظور کشف و استخراج دانش مورد تحلیل و کندوکاوهای خودکار و نیمه خودکار قرار می‌دهد. با استقرار پایگاه داده در سازمان مربوطه، زیرساخت‌های امکان داده‌کاوی با ابزارهای مربوطه برای سازمان پیدا خواهد شد. امروزه استفاده از رایانه‌ها در تحلیل و ذخیره‌سازی داده‌ها نقش بسزایی دارند. پایگاه‌های داده شامل چند صد میلیارد رکورد ثبت شده وجود دارند که امکان تحلیل و استخراج اطلاعات با روش‌های معمول و کلاسیک آماری از این پایگاه داده‌ها مستلزم داشتن دانش و ابزارهای توانمندتر است. شدت رقابت در عرصه‌های علمی، اجتماعی، اقتصادی، سیاسی و نظامی نیز اهمیت سرعت یا زمان دسترسی به اطلاعات را دو چندان کرده است. نیاز به طراحی سیستمهایی که قادر به اکتشاف سریع اطلاعات بدون خطای انسانی باشد مورد علاقه مدیران عصر حاضر است. در حال حاضر، داده‌کاوی مهمترین فناوری برای بهره‌برداری موثر، صحیح و سریع از داده‌های حجیم بوده و اهمیت آن رو به فزونی است و به‌عنوان ابزاری توانمند نه تنها دسترسی به اطلاعات را تسهیل می‌سازد بلکه باعث می‌شود تا الگوها و مدل‌های پنهانی در دل این پایگاه‌های داده نمایان شود و اطلاعات مفید و قابل اعتمادی که تا کنون نهفته بوده را به آشکار سازد.

⁴ Official Statistics

۱-۳- اهداف

هدف از این طرح پژوهشی معرفی روش‌های داده‌کاوی و به کارگیری برخی از این روش‌ها در آمارهای رسمی کشور است. پس از معرفی روش‌های داده‌کاوی، داده‌های هزینه و درآمد خانوار به‌طور خاص مورد مطالعه قرار می‌گیرد. خانوارهای شهری در سال ۱۳۹۰ و در سال ۱۳۹۶ به خوشه‌های ده‌گانه تقسیم‌بندی می‌شوند و رفتار خوشه‌ها با استفاده از قوانین داده‌کاوی کشف می‌گردد و در نهایت وابستگی این خوشه‌ها بررسی می‌گردد.

۱-۴- پیشینه تحقیق

داده‌کاوی دانشی بین‌رشته‌ای است و ترکیبی از علوم مانند آمار، ریاضیات، کامپیوتر، علوم اطلاعات، نظریه پایگاه داده و یادگیری ماشین می‌باشد. در قرن هجدهم قضیه بیز توسط توماس بیز چاپ و منتشر گردید که بعدها احتمال شرطی نام گرفت. قرن بعد قرن پیدایش الگوریتم رگرسیون توسط گاوس بود که به کمک آن پیش‌بینی‌هایی از داده‌های خام برای امور آتی انجام می‌شد. در قرن بعد الگوریتم شبکه‌های عصبی توسط دانشمندان کشف و منتشر شد. این ایده برگرفته از شبکه‌های عصبی مغز بود. برای اولین بار در سال ۱۹۸۰ اصطلاح «داده‌کاوی» مطرح شد. در این دوره کارشناسان می‌توانستند به روابط معنادار پی ببرند. سپس، در سال ۱۹۸۹ اصطلاح «کشف دانش در پایگاه داده» مطرح شد. در همین زمان اولین کارگاه آموزشی با نام KDD برای استخراج دانش در پایگاه‌های داده گوناگون، شروع به فعالیت نمود. در سال ۲۰۰۱ علم داده به عنوان یک رشته مستقل معرفی شد.^۵ امروزه داده‌کاوی در اقتصاد، مهندسی، سیاست، پزشکی و... کاربرد فراوانی دارد. داده‌کاوی تراکنش‌های مالی، سود سهام، شبکه‌های عصبی، امنیت ملی، الگوریتم‌های ژنتیک نمونه‌هایی از این کاربردهاست. در سال ۲۰۱۴ آقای حسنی و همکاران به بررسی اهمیت و ارزش کاربردهای داده‌کاوی در آمار رسمی پرداختند [۳]. در این طرح پژوهشی رفتار مصرفی خانوارها در سال‌های ۱۳۹۰ و ۱۳۹۶ بررسی و تغییرات آن با استفاده از تکنیک‌های داده‌کاوی نمایان خواهد شد.

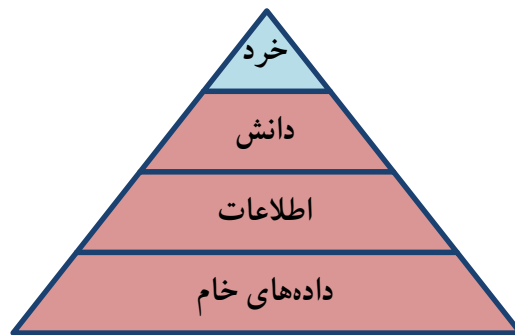
۱-۵- روش تحقیق

آنچه امروزه اهمیت بسیاری پیدا کرده کمبود یا نبود داده‌های مورد نیاز تحقیق نیست بلکه کمبود یا نبود روش‌هایی مناسب و استاندارد به منظور نگهداری، به روز کردن، در دسترس قرار دادن و در حالت آرمانی‌تر، کشف دانش جدید در داده‌های موجود است. یکی از راهکارهای پیشنهادی برای حصول به این هدف، استفاده از روش‌های داده‌کاوی است. روش‌های داده‌کاوی این امکان را به کاربر می‌دهند که بتواند انبوه داده‌های جمع‌آوری شده را تحلیل و تفسیر نماید و الگوها و دانش نهفته در آن را استخراج کند. از آن‌جا که مراکز و مؤسسه‌های آماری با انبارداده‌های با حجم بالا سروکار دارند، خواسته یا ناخواسته نیازمند استفاده از این دانش سودمند می‌باشند. در حال حاضر، داده‌کاوی مهمترین فناوری برای بهره‌برداری

⁵ Knowledge Discovery in Databases

⁶ William S. Cleveland

موثر، صحیح و سریع از داده‌های حجیم بوده و اهمیت آن رو به فزونی است و به عنوان ابزاری توانمند نه تنها دسترسی به اطلاعات را تسهیل می‌سازد بلکه باعث می‌شود تا الگوها و مدل‌های پنهانی در دل این پایگاه‌های داده نمایان شود و اطلاعات مفید و قابل اعتمادی که تا کنون نهفته بوده را به آشکار سازد. همان‌طور که پیش‌تر اشاره شد، امروزه داده‌کاوی در دو بخش دولتی و بخش خصوصی نقش حیاتی دارد. در بخش خصوصی (بانکداری، بیمه، بهداشت، بازاریابی و...) برای کاهش هزینه‌ها، ارتقاء کیفی پژوهش‌ها و بالاتر بردن میزان فروش از داده‌کاوی استفاده می‌نمایند. در بخش دولتی نیز به منظور برنامه‌ریزی و تعیین سیاست‌های کلان برای اداره جامعه و آن سازمان دولتی از داده‌کاوی استفاده می‌شود. ابزارهای داده‌کاوی ممکن است شامل مدل‌های آماری الگوریتم‌های ریاضی و روش‌های یادگیرنده باشد. داده‌کاوی منحصر به گردآوری و مدیریت داده‌ها نبوده و تجزیه و تحلیل اطلاعات و پیش‌بینی را نیز شامل می‌شود. داده‌کاوی با سه قسمت اولیه‌ی هرم دانش ارتباط تنگاتنگی دارد. این بخش‌ها در شکل ۱ به رنگ نارنجی نشان داده شده است.



شکل ۱-۱- هرم دانش

به کمک داده‌کاوی هم می‌توان روش‌های توصیفی مانند خوشه‌بندی، کشف قواعد پیوند و کشف الگوهای ترتیبی را انجام داد و هم می‌توان روش‌های پیش‌بینی کننده مانند رگرسیون، رده‌بندی را انجام داد. برخی از روش‌های مهم داده‌کاوی عبارتند از:

- رده‌بندی (Classification): رده‌بندی یعنی هر متغیری در یک کلاس خاص قرار گیرد. هر کلاس شامل مجموعه‌ای از متغیرهاست. یعنی نیازمند مدلی هستیم که کلاس متغیرها را به عنوان خروجی یک تابع توصیف نماید. این الگوریتم تحت نظارت کاربر صورت می‌پذیرد. مهم‌ترین الگوریتم‌های رده‌بندی، درخت تصمیم، شبکه عصبی و روش بیز می‌باشند.
- خوشه‌بندی (Clustering): گروه‌بندی مجموعه‌ای از اعضاء، رکوردها یا اشیاء به نحوی که اعضای موجود در یک خوشه بیشترین شباهت را به یکدیگر و کمترین شباهت را به اعضای خوشه‌های دیگر داشته باشند. این الگوریتم تا حدی بدون نظارت کاربر صورت می‌پذیرد.
- قواعد پیوند (Association Rules): قواعد انجمنی که به تحلیل سبد بازار نیز معروف است. الگوریتمی است که بر اساس آن وابستگی رویدادی به رویداد دیگر نمایش داده خواهد شد. به‌عنوان مثال خرید کره و مربا به هم وابسته‌اند.

- رگرسیون (Regression): پیش‌بینی یک متغیر براساس سایر متغیرها بر مبنای یک مدل وابستگی خطی یا غیرخطی، رگرسیون نام دارد.
- تحلیل‌های دنباله‌ی (Sequence Analysis): تجزیه و تحلیل رویدادهایی است که به صورت سلسه مراتبی رخ می‌دهد.

در واقع الگوریتم‌های داده‌کاوی مبتنی بر دو اصل است:

- الگوریتم‌های توصیفی: یافتن الگوها و روابط بین متغیرها، روش‌های مرتبط با آن.
 - الگوریتم‌های پیش‌بینی: یافتن مدل‌های مناسب و نتایجی مفید برای تصمیم‌گیری‌های آتی.
- در واقع به کمک داده‌کاوی هم می‌توان روش‌های توصیفی مانند خوشه‌بندی، کشف قواعد پیوند و کشف الگوهای ترتیبی را انجام داد و هم می‌توان روش‌های پیش‌بینی‌کننده مانند رگرسیون، دسته‌بندی را انجام داد. با توجه به ضرورت به‌کارگیری داده‌کاوی در آمار رسمی، در این طرح به طور خاص داده‌های هزینه و درآمد خانوار مورد مطالعه قرار می‌گیرد. خوشه‌های ده‌گانه با استفاده از قوانین داده‌کاوی کشف می‌گردد و وابستگی این خوشه‌ها نیز بررسی می‌شود.