

داده کاوی یک ابزار آنالیز مدیریتی

Data Mining an Analysis Implement Managemental

استاد راهنما : مهندس عمادی

ارائه دهنده :

مسعود فهرستی خمسه

یوسف مغانلو فرد

قربان مقدم زرزی

موسسه آموزش عالی روزبه زنجان

چکیده

داده کاوی فرآیندی تحلیلی برای کاوش داده های طراحی شده است، که در جستجوی الگوهای سازگار، یا روابط سیستماتیک بین متغیرها است، و سپس به تأیید این یافته ها با استفاده از الگوهای تشخیص داده شده می پردازد. استخراج اطلاعات مناسب از میان انبوه داده ها و تبدیل آنها به دانش مورد نیاز سازمانها، بویژه در تصمیم گیری های سازمانی، نیازمند استفاده از روش های نوین در این حوزه است. داده کاوی یکی از این ابزار و رویکردهاست که در فضای مدیریت دانش سازمانها به کشف دانش از پایگاه داده ها کمک می کند. این مقاله به بررسی ویژگی های منحصر به فرد این حوزه از فناوری و تکنیکهای استفاده از آن را نشان می دهد.

مقدمه

به مرور زمان، استخراج و کشف سریع و دقیق اطلاعات با ارزش و پنهان از پایگاه داده ها به عنوان داده کاوی مورد توجه قرار گرفت. به این شکل بود که فرایند داده کاوی به عنوان فرایند آماری و تجزیه و تحلیل در فرایند کشف دانش در پایگاه داده ها (KDD) پیرنگ شد، به حدی که گاه، داده کاوی (DM) به عنوان مترادف کشف دانش در پایگاه داده ها (KDD) مورد استفاده قرار می گرفت [2]. امروزه فرایند استخراج اطلاعات معتبر، از پیش ناشناخته، قابل فهم و قابل اعتماد از پایگاه داده های بزرگ و استفاده از آن در تصمیم گیری و در فعالیتهای تجاری داده کاوی نامیده می شود [1]. در تعاریف متعدد و متنوع برای داده کاوی بر موضوعاتی نظیر: استخراج دانش کلان، کاوش در داده ها، تجزیه و تحلیل داده ها و یافتن روابط و الگوهای مطمئن بین داده ها تأکید می شود. هدف نهایی داده کاوی، ایجاد سیستم های پشتیبانی تصمیم گیری سازمانی است. داده کاوی به استخراج اطلاعات مفید و دانش از حجم زیاد داده ها می پردازد.

داده‌کاوی، الگوهای حاوی اطلاعات را در داده‌های موجود جست‌وجو می‌کند. این الگوها و الگوریتم‌ها، می‌توانند توصیفی باشند یعنی داده‌ها را توصیف کنند و یا جنبه پیش‌بینی داشته باشند، یعنی از متغیرها برای پیش‌بینی ارزش‌های ناشناخته سایر متغیرها به کار روند. داده‌کاوی توصیفی، به دنبال یافتن اگرها در فعالیت‌ها یا اقدامات گذشته است و داده‌کاوی پیش‌بینانه با نگاه به سابقه، رفتار آینده را پیش‌بینی می‌کند [1].

حوزه فعالیت

اکتشاف در این مرحله معمولاً با آماده‌سازی داده‌ها که ممکن است شامل تمیز کردن داده‌ها، تبدیل داده‌ها، زیرمجموعه‌های انتخاب آثار ضبط شده و انجام برخی از عملیات اولیه انتخاب شروع می‌شود. سپس بسته به ماهیت تحلیلی، این مرحله از فرایند استخراج داده‌ها ممکن است شامل هر انتخاب ساده و سراسر برای یک مدل رگرسیون استادانه درست شده را به تجزیه و تحلیل اکتشافی با استفاده از طیف گسترده‌ای از روش‌های گرافیکی و آماری به منظور شناسایی متغیرهای مربوطه و تعیین پیچیدگی از طبیعت مدل‌ها باشد. البته ناگفته نماند که داده‌کاوی معمولاً با نوشتن مقدار زیادی گزارش و تحقیق و استعمال در آنها اشتباه گرفته می‌شود. اما در واقع داده‌کاوی هیچ‌کدام از اینها را شامل نمی‌شود. داده‌کاوی توسط تجهیزات خاصی صورت می‌پذیرد، که عملیات کاوش را بر اساس تجزیه و تحلیل مکرر داده‌ها انجام می‌دهد. داده‌کاوی با آنالیزهای متداول آماری نیز متفاوت است؛ در زیر به بررسی نوع دیدگاه روش آماری به روش داده‌کاوی می‌پردازیم:

روش آنالیز آماری :

یک مفسر ممکن است متوجه الگوی رفتاری شود که سبب کلاهبرداری بیمه‌گردد. بر اساس این فرضیه، مفسر به طرح یک سری سوال می‌پردازد تا این موضوع را بررسی کند. اگر نتایج حاصله مناسب نبود، مفسر فرضیه را اصلاح می‌کند و یا با انتخاب فرضیه دیگری مجدداً شروع می‌کند. این روش نه تنها وقت‌گیر است بلکه به قدرت تجزیه و تحلیل مفسر نیز بستگی دارد. مهمتر از همه اینکه این روش هیچ‌وقت الگوهای کلاهبرداری دیگری را که مفسر به آنها مظنون نشده و در فرضیه جا نداده، پیدا نمی‌کند.

روش داده‌کاوی :

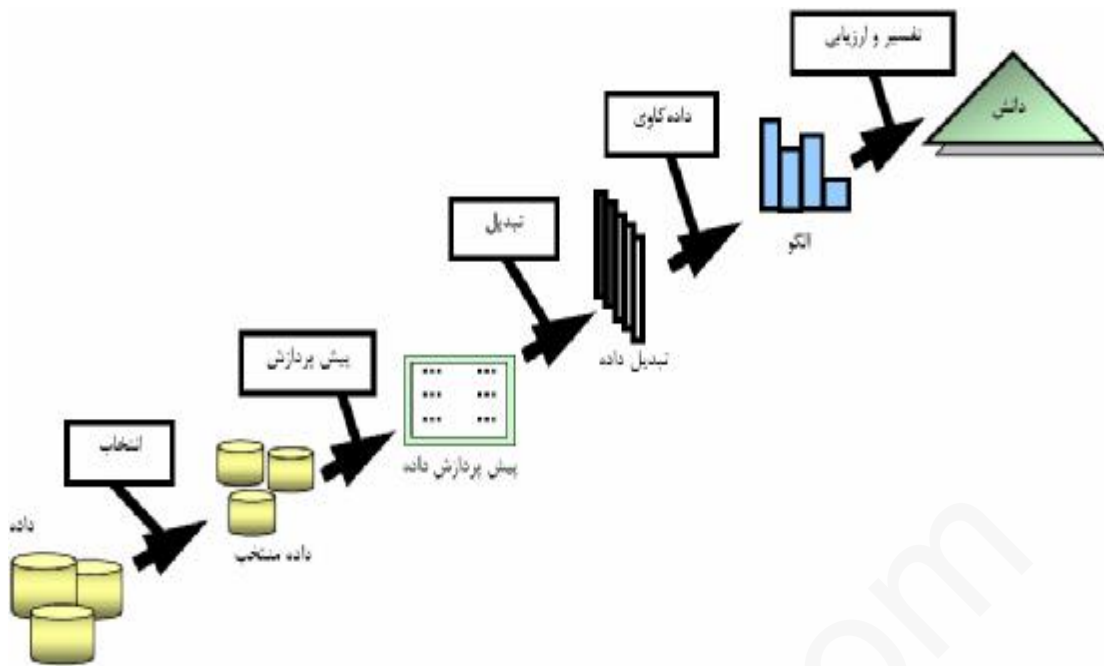
یک مفسر سیستم‌های داده‌کاوی را ساخته و پس از طی مراحل از جمله جمع‌آوری داده‌ها، یکپارچه‌سازی و اخلاص داده‌ها به انجام عملیات داده‌کاوی می‌پردازد. داده‌کاوی تمام الگوهای غیرعادی را که از حالت عادی و نرمال انحراف دارند و ممکن است منجر به کلاهبرداری شوند را پیدا می‌کند. نتایج داده‌کاوی حالت‌های مختلفی را که مفسر باید در

مراحل بعدی تحقیق کند، نشان می دهند. در نهایت مدل های به دست آمده می توانند مشتریانی را که امکان کلاهبرداری دارند، پیش بینی نمایند.

تحلیلهای داده کاوی به دو روش با ناظر و بدون ناظر و از طریق الگوریتمهایی چون شبکه های عصبی (NN)، طبقه بندی و درخت تصمیم (C&RT)، ژنتیک، تحلیل سبد خرید، شبکه کوهونن قابل اجراست. علاوه بر این الگوریتمهای رایج، همچنان الگوریتمهای جدیدی برای اهداف تحقیقات علمی یا تجاری از طریق طرحهای پژوهشی دانشگاهی، تولید می شود. ویژگیهای منحصر بفرد داده کاوی را می توان به صورت زیر برشمرد [3]:

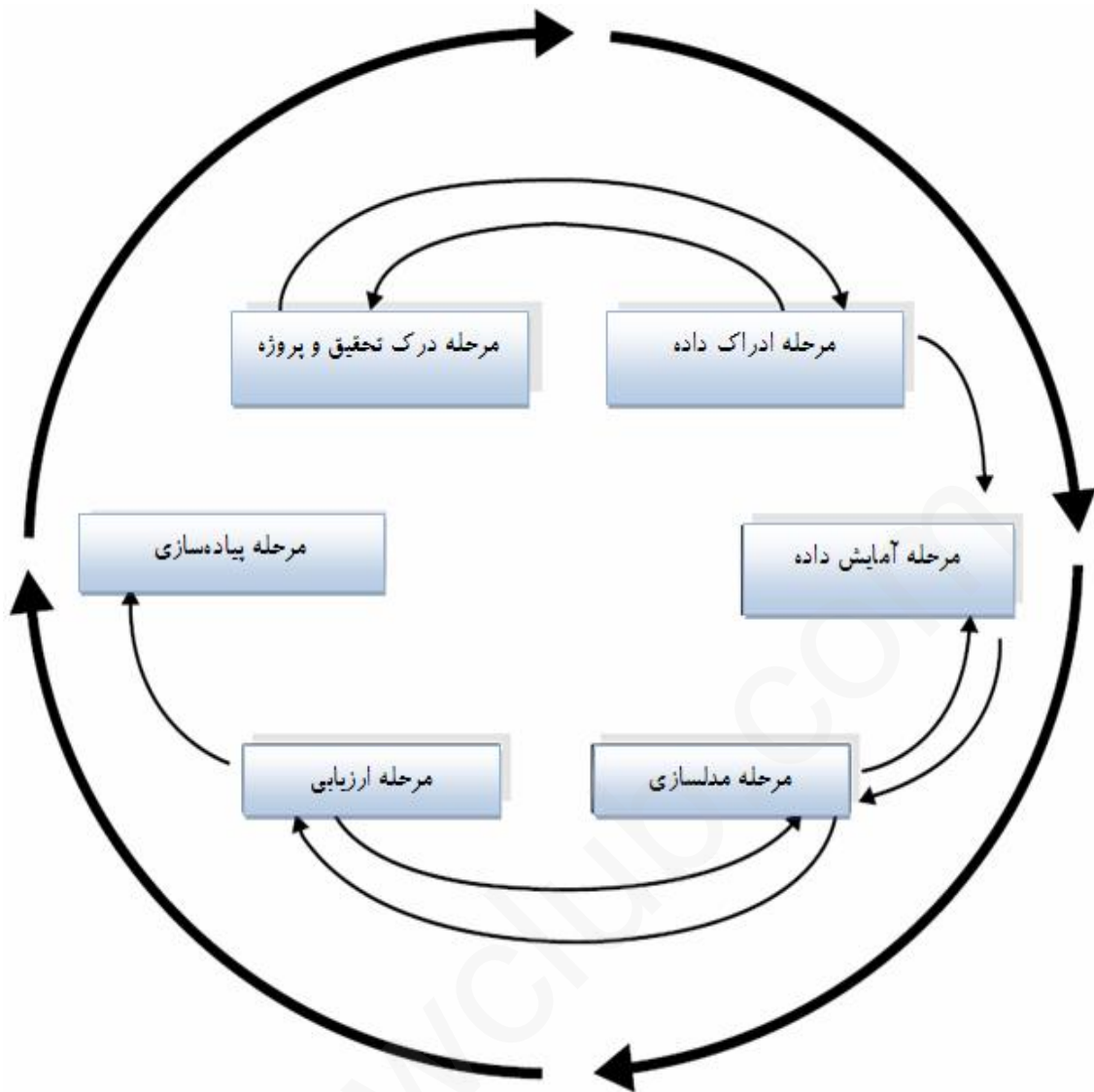
- نه تنها بر فاز تحلیل، بلکه بر طراحی مطالعه و جمع آوری داده نیز تاثیر می گذارند
- امکان جستجوی پاسخ سؤالات دقیق و با پیچیدگی بالا را در دادههای جمع آوری شده فراهم می کنند.
- قادرند که به سؤالات بطور واضح و مشخص پاسخ دهند. مزیت اصلی و تفاوت آنها با سایر تکنیکها نیز در همین است که بجای ارائه صرف استراتژی کلان، پاسخهای دقیق در اختیار محقق قرار می دهند.
- امکان سنجش اثر متغیرهای مختلف بر روی متغیرهای وابسته را فراهم می کنند.
- به مدیران کمک می کنند که تأثیر سناریوهای آتی را مورد ارزیابی قرار دهند و با مدلسازی گزینه های متعدد و کمک به تصمیم گیری در شرایط عدم قطعیت به انتخاب مسیر حرکت پردازند.

محققینی که تنها روابط دو به دو را در نظر می گیرند و از داده کاوی استفاده نمی کنند، ابزار قدرتمندی را از دست می دهند که می تواند اطلاعات سودمندی را در اختیار آنان قرار دهد. در مسائل واقعی چندین متغیر به طور همزمان بر روی پاسخ تاثیر می گذارند، از این رو آنالیزهای چندمتغیره جواب های دقیقتر و نزدیک به واقع تری را فراهم می کند. در شکل (۱) فرایند کسب دانش از پایگاه دادهها به صورت شماتیک بیان شده است [4] همانطور که ملاحظه می شود یکی از گام های این فرایند، داده کاوی می باشد. موفقیت در این مرحله کاملاً متأثر از سه گام قبل است بگونه ای که اگر هر کدام از مراحل قبلی به درستی انجام نپذیرد، نتایج حاصل از داده کاوی نه تنها مفید نبوده ممکن است گمراه کننده نیز باشد.



شکل (۱): فرآیند تبدیل داده‌ها به دانش

تکنیکهای داده کاوی از جمله تکنیکهای نوین علمی هستند که در توصیف، تشریح، پیش بینی و کنترل پدیدهها به کار می روند [3]. این تکنیکها به اندازه گیری، تشریح و پیش بینی درجه وابستگی میان متغیرها میپردازند. روشهای داده کاوی نه تنها بر جنبه های تحلیلی مطالعات، بلکه در طراحی و ابزارهای جمع آوری داده برای تصمیم گیری و حل مسائل نیز تأثیر می گذارند. موفق ترین پروژههای داده کاوی، در چارچوب فرآیند استاندارد اجرا می شود که توسط یک تیم کاری در شرکت SPSS در قالب پروژههای به نام CRISP-DM ارائه شده است [5]. برطبق CRISP-DM یک پروژه داده کاوی معین شامل چرخه حیات شش مرحله ای است که توالی مراحل را نشان می دهد شکل (۲). هر مرحله از ترتیب مراحل اغلب نتیجه وابستگی مراحل قبلی را نیز دربر دارد. مهمترین وابستگی بین مراحل نمایش پیکانها می باشد. خاصیت تکراری CRISP حاکی از چرخه بیرونی است که اغلب منجر به راه



شکل (۲) CRISP-DM در فرایند تکرار و سازگاری مراحل

حلی برای مسئله تحقیقی یا تجاری با سوالات اضافی جالب توجه می شود. در زیر مراحل کاری در داده کاوی را توضیح می دهیم:

مرحله درک پروژه و فهم حوزه کاربرد: اولین مرحله پردازش استاندارد CRISP-DM می باشد که به صورت آشکار اهداف و نیازمندیها آن مشخص می شود. ترجمه اهداف و محدودیت آن در قاعده سازی، تعریف مسئله داده کاوی و مهیا کردن استراتژی اولیه برای نائل شدن به اهداف تعریف می شود.

مرحله انتخاب داده ها: این مرحله شامل جمع آوری داده ها برای استفاده از تحلیل اکتشافی و مشخص کردن اطلاعات اولیه برای ارزیابی داده های با کیفیت و انتخاب داده های مفید و مورد نیاز می باشد.

مرحله آماده سازی داده‌ها: آماده کردن داده‌های اولیه خام به داده‌های نهایی، این داده‌ها در کلیه مراحل بعدی استفاده می‌شود و از این نظر این مرحله تحلیل و تلاش بیشتری را می‌طلبد. انتخاب عناصر و شناسه‌های تحلیل شده را برای کاوش داده‌ها اختصاص می‌دهیم. و با تمیز کردن داده‌های خام آن را برای ابزارهای مدلسازی آماده می‌کنیم.

مرحله مدلسازی: با انتخاب و به کار بستن تکنیکهای مدلسازی مناسب و روش داده‌کاوی معین نتایج مدلسازی را بهینه می‌کنیم که در صورت نیاز می‌توانیم با برگشت به عقب تحلیل مدلسازی را بهینه تر نماییم.

مرحله ارزیابی: مشخص کردن اینکه آیا مدل انتخابی، ما را به اهدافمان که در اولین مرحله تعیین کردیم می‌رساند. اتخاذ تصمیم راجع به استفاده از نتایج داده‌کاوی برای اعتبارسنجی نیز در این مرحله انجام می‌شود.

مرحله تحکیم و گسترش: استفاده کردن از مدل ایجاد شده، برای مثال می‌تواند تولید یک گزارش ساده از خروجیها را نام برد، و برای یک مثال پیچیده تکمیل کردن پردازش داده‌کاوی موازی در سایر حوزه‌ها می‌باشد که این الگوها به یک دانش مفید و قابل استفاده تبدیل می‌شوند و پس از بهبود آنها، الگوهایی که کارا محسوب می‌شوند در یک سیستم اجرایی به کار گرفته خواهند شد.

نتیجه گیری

بررسی اجمالی پژوهشهای صورت گرفته در حوزه دانش ابزارهای داده‌کاوی نشان می‌دهد که تحقیقات عمیق و اساسی در این باره خصوصاً در ایران اندک شمار است. از سوی دیگر با افزایش سرعت تحول در علوم، ضرورت استفاده از دانشهای نوین بیش از پیش محرز شده است. داده کاوی به عنوان یک رشته علمی نوین در زمینه بازاریابی و استخراج اطلاعات می‌تواند نقش مهمی در جهت دستیابی به این اهداف داشته باشد. امروزه اکثر نرم افزار های پایگاه داده ای مثل ORACLE و SQL Server نیز شامل ابزارهایی داده کاوی شده اند ولی نرم افزار های تخصصی داده‌کاوی همچون Knowledge Intelligent Miner , Darwin , Mine Set Studio, Data Mind از مهمترین ابزار های داده کاوی به شمار می‌روند. در این مقاله قابلیت‌های داده کاوی و مراحل کاری آن معرفی شد که در گامهای بعد می‌توان تأثیر آن را در عمل آزمود.

مراجع

- [1] B. Fernandez / Et. Al., "Knowledge Management"/ Cho. 12, 2004.
- [2] N.Balac/ "Introduction To Data Mining" , 2006
- [3] Hair ,Joseph F., "Multivariate Data Analysis", Prentice Hall, 2005.

- [4] Daniel T. Larose, *"Discovering Knowledge in Data: An Introduction to Data Mining"*, 2004 .
- [5] www.spss.com/ CRISP DM/ Downloads
- [6] Pang-Ning Tan, Steinbach, *"Introduction to Data Mining"*, 2005 .

knowclub.com