

معرفی داده کاوی یا Data mining

معرفی داده کاوی یا data mining داده کاوی یا کشف دانش در پایگاه داده ها (KDD) علم نسبتاً تازه ای است که با توجه پیشرفت کشور در زمینه IT و نگاه های ویژه به دولت الکترونیک و نفوذ استفاده از سیستم های رایانه ای در صنعت و ایجاد بانک های اطلاعاتی بزرگ توسط ادارات دولتی، بانک ها و بخش خصوصی نیاز به استفاده از آن به طور عمیقی احساس می شود .

معرفی داده کاوی

داده کاوی یا کشف دانش در پایگاه داده ها (KDD) علم نسبتاً تازه ای است که با توجه پیشرفت کشور در زمینه IT و نگاه های ویژه به دولت الکترونیک و نفوذ استفاده از سیستم های رایانه ای در صنعت و ایجاد بانک های اطلاعاتی بزرگ توسط ادارات دولتی، بانک ها و بخش خصوصی نیاز به استفاده از آن به طور عمیقی احساس می شود. داده کاوی یعنی کشف دانش و اطلاعات معتبر پنهان در پایگاه های داده. یا به بیان بهتر تجزیه و تحلیل ماشینی داده ها برای پیدا کردن الگوهای مفید و تازه و قابل استناد در پایگاه داده های بزرگ ، داده کاوی نامیده می شود. داده کاوی در پایگاه های داده کوچک نیز بسیار پرکاربرد است و از نتایج و الگوهای تولید شده بوسیله آن در تصمیم گیری های استراتژیک تجاری شرکتهای کوچک نیز می توان بهره های فراوان برد. کاربرد داده کاوی در یک جمله را این گونه می توان بیان کرد : " داده کاوی اطلاعاتی می دهد ، که شما برای گرفتن تصمیم هوشمندانه ای درباره مشکلات سخت شغلستان به آنها نیاز دارید. "

مثالی کلاسیک از کاربرد داده کاوی

اغلب تجارت ها به تصمیم گیریهایی استراتژیک و یا اتخاذ خط مشی های جدید برای خدمت رسانی بهتر به مشتریان نیاز دارند. به عنوان مثال فروشگاهها آرایش مغازه خود را برای ایجاد میل بیشتر به خرید مجدداً طراحی می کنند. این مثال به داده هایی در مورد رفتار مصرفی گذشته مشتریان برای تعیین الگوهایی به وسیله داده کاوی، نیاز دارند.

برای روشن تر شدن مسئله می توان مثال را اینگونه بیان کرد که در یک فروشگاه زنجیره ای پس از داده کاوی مشخص میشود که درصدی از مشتریان خرید تلوزیون ، میز تلوزیون و گلدان کریستالی را هم در همان روز و بعد از خرید تلوزیون میخرند.مدیر فروشگاه می تواند بلافاصله دستوراتی صادر کند که براساس مدلهای تلوزیون موجود میزهایی و براساس مدل میزها گلدانهای کریستالی برای فروش سفارش داده شود و غرفه های جنبی غرفه تلوزیون را به میز و گلدان کریستالی اختصاص دهد. مطمئناً حتی پس از مدت کوتاهی سود حاصل از این بخش از فروشگاه به طور قابل ملاحظه ای ترقی خواهد کرد.

در واقع ابزار داده کاوی، داده را می گیرد و يك تصویر از واقعیت به شکل مدل می سازد، این مدل روابط موجود در داده ها را شرح می دهد .

برای بهبود بهره وری از یک فروشگاه داده کاوی از داده های انبار داده ، مدل هایی را ارائه میدهد که بیانگر این هستند که چه محصولات یا خدماتی، به چه مشتریانی، در چه زمانی و از طریق چه کانالی عرضه شود .

بیشتر شرکتهای، بانکهای داده ای عظیمی شامل داده های بازاریابی، منابع انسانی و مالی را دارا هستند. بنابراین، سرمایه گذاری در زمینه انبار داده، یکی از اجزای حیاتی در استراتژی مدیریت ارتباط با مشتری است .

رابطه مشتری با زمان تغییر می کند و چنانچه تجارت و مشتری درباره یکدیگر بیشتر بدانند این رابطه تکامل و رشد می یابد. چرخه زندگی مشتری چارچوب خوبی برای به کارگیری داده کاوی در مدیریت ارتباط با مشتری فراهم می کند. در بخش ورودی داده کاوی، چرخه زندگی مشتری می گوید چه

اطلاعاتی در دسترس است و در بخش خروجی آن، چرخه زندگی می گوید چه چیزی احتمالاً جالب توجه است و چه تصمیماتی باید گرفته شود. داده کاوی می تواند سودآوری مشتری های بالقوه را که می توانند به مشتریان بالفعل تبدیل شوند، پیش بینی کند و اینکه تا چه مدت به صورت مشتریان وفادار خواهند ماند و چگونه احتمالاً ما را ترک خواهند کرد.

بعضی از مشتریان مرتباً مراجعاتشان را به شرکتها برای کسب مزیتهایی که طی رقابت میان آنها به وجود می آید، تغییر می دهند. در این صورت شرکتها می توانند هدفشان را روی مشتریانی متمرکز کنند که سودآوری بیشتری دارند. بنابراین می توان از طریق داده کاوی ارزش مشتریان را تعیین، رفتار آینده آنها را پیش بینی و تصمیمات آگاهانه ای را در این رابطه اتخاذ کرد.

از کاربرد های داده کاوی می توان به نمونه های زیر اشاره کرد:

1. بانکداری:

- از جالب توجه ترین کاربردهای داده کاوی می توان به کشف پول شویی اشاره کرد .
- تشخیص مشتریان ثابت و همیشگی
- تعیین مشتریان استفاده کننده از یک سرویس خاص

2. بیمه:

- پیش گویی میزان استقبال از بیمه نامه های جدید
- تشخیص کلاهبرداری ها و مشخص کردن رفتار های نامتناسب
- تشخیص نیاز مشتریان و خواسته های آنها
- تشخیص تخلفات پزشکی

واضح است که زمینه استفاده از داده کاوی بی نهایت گسترده است. و دو مثال فوق به خاطر درک راحت تر انتخاب شده اند.

داده کاوی شباهت زیادی به تحلیل های آماری دارد. ولی داده کاوی از جهات زیادی با آمار متفاوت است و مزیت های زیادی نسبت به آمار دارد. جالب ترین تفاوت داده کاوی با تحلیل های آماری این است که در آمار ما فرضیه ای طرح می کنیم و با استفاده از تحلیل های آماری به اثبات یا رد فرضیه می پردازیم اما داده کاوی به فرضیه احتیاجی ندارد. در واقع ابزار داده کاوی فرض می کند که شما خود هم نمی دانید به دنبال چه می گردید. و این نکته ای است که باعث می شود کار آمدی داده کاوی در مواقع بروز مشکل نمایان شود . برای مثال ما در آمار فرض می کنیم که دو گروه فاصله ای باهم ارتباط دارند سپس با استفاده از ضریب هم بستگی پیرسون مشخص می کنیم که ارتباط وجود دارد یا خیر . ولی داده کاوی بدون توجه به اینکه ما اینگونه فرضی داشته باشیم یا نه با کاوش میان داده ها اگر ارتباطی مخفی معنی داری وجود داشته باشد آن را به اطلاع ما می رساند . تفاوت بعدی آمار و داده کاوی در این است که آمار فقط می تواند از داده های عددی استفاده کند ولی داده کاوی از داده های غیر عددی هم استفاده می کند . تفاوت های دیگری هم میان آمار و داده کاوی وجود دارد که بحث در مورد آنها در حوصله این مقاله نمی گنجد.

اما برای اولین بار در سال ۱۹۵۰ از رایانه برای تحلیل و ذخیره پایگاه داده ها استفاده شد. ولی حجم اطلاعات و میزان رشد آنها به قدری زیاد بوده است که هم اکنون کسی از میزان اطلاعات ذخیره شده در پایگاه داده های سراسر دنیا به صورت دقیق اطلاعی ندارد ولی مطمئناً حجم اطلاعات و مخصوصاً

سرعت رشد آنها به قدری زیاد شده که آمار شناسان و تحلیل گران در بررسی و تحلیل پایگاههای داده در زمینه های مختلف ناتوانند. بعضی از پایگاه داده ها به قدری بزرگ و پیچیده شده اند که تحلیل روابط و استخراج اطلاعات مفید پنهان شده در آنها واقعا از ظرفیت ذهنی بشری فراتر رفته است. از زمانی که رشد پایگاه های داده و حجم اطلاعات، سرعت گرفت و میزان داده ها افزایش یافت، نیاز به تحلیل ماشینی داده ها و استخراج سریع و دقیق دانش نهفته در آنها احساس شد. شاید بتوان لوول (۱۹۸۳) را اولین شخصی دانست که گزارشی در مورد داده کاوی تحت عنوان « شبیه سازی فعالیت داده کاوی » ارائه نمود.

عمل داده کاوی از یک پایگاه داده به چند مرحله مشخص تقسیم می شود که ما در این مقاله به معرفی و توضیحی مختصر در مورد هر یک از این مراحل اکتفا می کنیم:

1. مرحله اول : تشکیل انبار داده.

با توجه به عنوان، این مرحله برای تشکیل محیطی پیوسته و یک پارچه جهت انجام مراحل بعدی و داده کاوی در آن، انجام می گیرد. در حالت کلی انبار داده مجموعه پیوسته و طبقه بندی شده است که دائما در حال تغییر بوده و دینامیک است که برای کاوش آماده می شود.

2. مرحله دوم : انتخاب داده ها.

در این مرحله برای کم کردن هزینه های عملیات داده کاوی، داده هایی از پایگاه داده انتخاب می شوند که مورد مطالعه هستند و هدف داده کاوی دادن نتایجی در مورد آنهاست.

3. مرحله سوم : تبدیل داده ها.

مشخص است برای انجام عملیات داده کاوی لزوما باید تبدیلات خاصی روی داده ها انجام گیرد ممکن است این تبدیلات خیلی راحت و مختصر مثل تبدیل byte به integer باشد یا خیلی پیچیده و زمان بر و با هزینه های بالا مثل تعریف صفات جدید و یا تبدیل و استخراج داده ها از مقادیر رشته ای و ... باشد.

4. مرحله چهارم : کاوش در داده ها.

در این مرحله است که داده کاوی انجام می شود. در این مرحله با استفاده از تکنیک های داده کاوی داده ها مورد کاوش قرار گرفته، دانش نهفته در آنها استخراج شده و الگو سازی صورت می گیرد.

5. مرحله پنجم : تفسیر نتیجه.

در این مرحله نتایج و الگو های ارائه شده توسط ابزار داده کاو مورد بررسی قرار گرفته و نتایج مفید معین می شود.

طرز کار ابزار داده کاو اینگونه است که ابزار به دنبال اثبات این است که وجود چیزی به معنای وجود چیز دیگری است و سعی می کند در درجه اول از توالی ارتباطات برای کشف یک الگو بهره بگیرد و در نهایت اطلاعات بدست آمده را دسته بندی کند تا به الگوی خاصی برسد که بتواند آن را براساس فاکتورهای داخلی به مخاطبش ارائه دهد.

همچنین در داده کاوی از الگوریتم های ژنتیک و شبکه های عصبی هم استفاده می شود. شبکه های عصبی به علت کار آمدی در حل مسائل پیچیده و بزرگ مورد استفاده اند و کاربرد الگوریتم های ژنتیک در داده کاوی برای جستجو و ساختن یک مدل بهینه در میان مدل های بدست آمده است، به این گونه که مدل های اولیه روی کروموزوم هایی قرار می گیرند و با رقابت بر سر انتقال صفات به نسل بعد، بهترین مدل و لایق ترین آنها به کاربر ارائه می شوند.

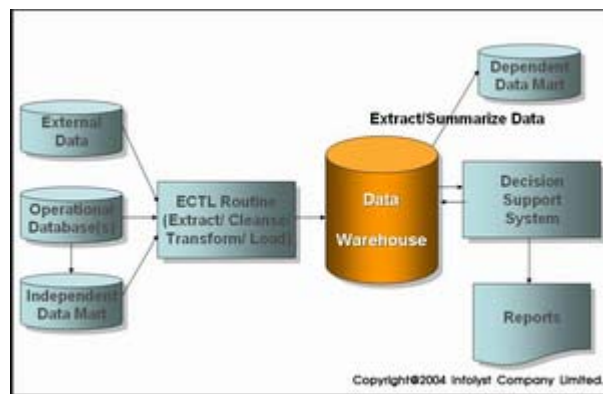
داده کاوی امروز گسترش زیادی یافته است به طوری که اکثر نرم افزار های پایگاه داده ای مثل SQL Server و ORACLE نیز شامل ابزارهایی داده کاوی شده اند ولی هنوز نرم افزار های تخصصی داده کاوی همچون Intelligent Miner, Darwin, Mine Set, Knowledge Studio, Data Mind از مهمترین ابزار های داده کاوی اند.

آشنایی با مفاهیم انبار های داده (Datawarehouse)

کاوی داده / انبار داری / تحلیل زنده

مقدمه

آوری و فن (Datawarehouse) **داده** مفاهیم انبار های **یا** آشنایی خواننده **مقاله** هدف از ارایه این برای ساخت و بهره برداری از برنامه پایه گذار فن آوری لازم **داده** های مرتبط می باشد. انبار های تخیلی با آنها آشنا شده ایم. برنامه هایی که های هوشمندی هستند که بعضا در برخی فیلم های نموده و پس از تجزیه و تحلیل با آنها به مشاوره می پردازند اطلاعات را برای صاحبان خود جمع آوری



چيست؟ OLTP

در دهه فقط در سال ۲۰۰۰ میزان ظرفیت نصب شده جهت ذخیره سازی اطلاعات از کل ظرفیت موجود ۱۹۹۰. بیشتر بوده است

تقریباً حجم کل اطلاعات در کامپیوترها هر ۵ هاست. در حال حاضر **داده** حیات بازرگانی نوین مبتنی بر ایجاد برنامه های چند رسانه ای و بانکهای اطلاعاتی پیش سال دو برابر می شود و با توجه به سرعت اطلاعات به دو برابر در سال برسد بینی می شود که شتاب رشد جدیدی هستند که امور خود را توسط کامپیوترها تولید کنندگان این اطلاعات موسسات و شرکت های جمع آوری نموده و به مصرف می ها را **داده** مکانیزه ای که هدایت می کنند. سیستم های تولید ها هستند **داده** تولید کنندگان واقعی نامیده می شوند. این سیستم ها OLTP رسانند سیستم های

برنامه های کاربردی خادم و مخدوم بدو دسته تقسیم می شوند
(DSS) پشتیبانی تصمیم گیری سیستم های "
اطلاعات (OnLine) سیستم های پردازش زنده "

قبل از آنکه به این دو دسته هر یک راه های کاملاً متفاوتی را جهت حل مسایل تجاری ارایه می کنند. تفاوتهای این دو را بشناسیم بی ببریم لازم است **داده** ارزش انبار های

رزرواسیون، در کلیه خدمات بازرگانی دیده می شوند از جمله سیستم های OLTP سیستم های به زمان پاسخی بین ۱ تا دستگاه های فروش، کنترل انبار، سهام و فروش و ... این سیستم ها غالباً در ساعات مختلف روز، هفته و ماه می تواند ۳ ثانیه در ۱۰۰ در صد اوقات نیاز دارند. تعداد کاربران آنها همان زمان پاسخ قبلی نیاز دارند. در این گونه سیستم ها بشدت متغیر باشد و در تمامی این اوقات به خادمین تعاملی (Database Servers) با بانکهای اطلاعاتی معمولاً مخدومین بجای ارتباط متصل می شوند. البته این گونه ارتباط لازمه دستیابی به سرعت مورد نیاز (Transaction Servers) است. (Clients) مخدومین

قادرند تعامل را تقسیم می‌گردد. خادمین عادی (Heavy) و قوی (Lite) خود نیز به دو نوع عادی OLTP به اجرا بگذارند و خادمین قوی (StoredProcedures) در غالب پردازش های ثبت شده در بانک اطلاعاتی سربرابر برای دستیابی به سرعت OLTP میکنند. در برای اجرای دستورات استفاده (TP Monitor) از انتقال یک دستور ارتباطی شبکه ها در حداقل ممکن نگاه داشته می شود و غالبا ارتباطات در حد (سیکویل هستند SQL))

بسرعت پایگاه های اطلاعاتی بزرگی یا با جمع آوری اطلاعات امروزه حتی کوچکترین تجارتها هم قادرند چه رسد به وب سرور ها که می توانند ظرف مدت بسیار کوتاهی چندین صندوق های فروش ایجاد کنند. جمع آوری نمایند گیگا بایت اطلاعات

امروزه هر کسی بسادگی زمانی برای هر کار مکانیزه ای نیاز به میلیونها پول و ده ها متخصص بود. اما از امکانات رایانه ای بهره مند گردد. با خرید چند کامپیوتر شخصی و استخدام یک برنامه نویس می تواند ها برای همگان داده ایجاد پایگاه های خصوصی از عبارت دیگر دسترسی به خدمات رایانه ای برای شده است آسانتر

جمع آوری می شود مستقیما مورد استفاده OLTP توسط سیستم های هابی که داده در مجموع دانند چگونه نیاز ها چیستند و همچنین می داده قرار دارد. آنها دقیقا می دانند این افراد ایجاد کننده آن کنند های اطلاعاتی لحظه ای خود را که بطور روزمره بوجود می آید حل

داشته باشد به این اطلاعات نیاز OLTP سوالی که مطرح است اینست که اگر کسی خارج از مجموعه بایستی آنرا پیدا کرد و چگونه به ای موجود است؟ کجا داده چه باید کرد. این افراد از کجا می دانند چه است؟ چه معنایی دارد؟ آخرین چیزی که (Format) ها به چه شکلی داده آن دسترسی پیدا کنند؟ اجازه دهند دیگران به اطلاعات گرانهای آنان دسترسی به آن رضایت خواهند داد آنست که OLTP افراد دانند چه می خواهند، درخواستهای سیکویل زمانگیری را بر روی داشته باشند. کسانی که حتی نمی آورد ها را پایین می داده می کنند که سرعت و قابلیت سیستم تولید کننده بانکهای اطلاعاتی اجرا

مشابه خود در خود می خواستند با همکاران MIS در گذشته افراد بیرون از سیستم، از همکاران استخراج نمایند. اما امروزه حتی سیستم مربوطه تعامل داشته و نهایتا اطلاعات مورد نظر را از سیستم سازمان موجود است. اطلاعات بشدت توزیع خود هم بدرستی نمی داند چه اطلاعاتی در MIS مجموعه بخشی از اطلاعات سازمان وجود دارد شده و پراکنده است و تقریبا روی هر کامپیوتری

افراد یکی از ویژگیهای کامپیوتر های شخصی و همچنین معماری خادم/مخدوم موجب شده است که ترجیح می دهند اطلاعات اکثرا به اطلاعات سازمانی و کاربرد اطلاعات در سازمان علاقه ای نداشته و سازمان و شخصی (یا واحد های را تحت مالکیت شخصی اداره کنند به این ترتیب بین اطلاعات استخراج های سیستمی و اطلاعات داده متشکله) شکاف وجود خواهد داشت. از طرف دیگر بین می کنند افرادی شده نیز شکاف دیگری مشاهده می شود. کسانی که از بیرون به این اطلاعات نگاه بتوانند تصمیمات بهتری ها هستند بطوریکه داده هستند که بدنال یافتن طرحها، روالها و تمایلات در تجارت دیگران است و خیلی زود همگان بگیرند. تنیدن حصار بدور اطلاعات بمعنی تنیدن حصار در برابر بازنده جنگ این حصارها خواهند بود

چگونه اطلاعات را در اختیار داریم اگر بديگران اجازه دسترسی به آنرا بدهیم :اند شوند و از آن جمله داده سولات زیادی مطرح هستند که بایستی پاسخ نمی کند؟ چگونه مطمئن شویم که عملکرد بیرونی ها (غریبه ها) عملکرد سیستم ما را کند چه اطلاعاتی را بایستی در اختیار بیرونی ها قرار دهیم؟ است؟ (داده درونی و شخصی) فقط مربوط به سیستم تولید کننده چه اطلاعاتی اطلاعات به اشتراک گذاشته شده است؟ چه کسی مالک چه کسی این اطلاعات را بروز میکند؟ کنیم؟ بایستی بگذاریم دسترسی به اطلاعات مستقیم باشد یا آنرا در بانک دیگری کپی آیا اطلاعات استخراج شده چگونه نگهداری شده و چگونه بروز می شود؟

بشناسیم و تفاوتهای برای پاسخ به سوالات فوق بایستی نیاز های استفاده کنندگان از این اطلاعات را را درک کنیم OLTP میان سیستمهای پشتیبان تصمیم گیری و

میکند؟ ها استفاده داده چه کسانی از این

اطلاعات هستند (کسانی بیایید نامی برای این دسته از افراد انتخاب کنیم. این افراد مصرف کنندگان شکارچی اطلاعات می گذاریم چون هستند که تصمیمات استراتژیک می گیرند) فعلا نام این افراد را

نیازمند اطلاعات است. البته بازرگانان و دسترسی دارد و PC این نام معرف هر کسیست که به یک هستند صنعتگران اولین دسته از این افراد

سیستم پشتیبانی تصمیم گیری چیست؟

ها، تولید گزارش های **داده** ها ، یافتن ارتباط بین **داده** یک سیستم کارآمد، ابزار است برای تحلیل پاسخ به **قابلیت** اطلاعات در انواع ممکن، ها، راهکار های نمایش **داده** کارآمد، دسترسی منعطف به . ها به صفحات گسترده **داده** سوالات اگر ... چه ، چاپ اطلاعات، انتقال بیشتری در زمان پاسخگویی برخوردار **داده**، این ابزارها از انعطاف در مقایسه با سیستم های تولید است و قابلیت دسترسی همزمان کاربران به آن هستند. معمولا کنترل یکپارچگی در آنها رعایت نشده بروز رسانی اطلاعات غالبا بمعنی پردازش روی تمامی غالبا محدود است. جستجوی اطلاعات و یا برای غیر برنامه نویسان تهیه شده و بیشتر فعالیت ها در آن از طریق اطلاعات خواهد بود. این برنامه ها انجام می شود (Point and Click) کن نشان بده و کلیک

سیستم های اطلاعات مدیران اجرایی (Executive Information Systems)

همچنین به یک زمینه . قوی تر، ساده تر و کارآمدتر هستند DSS این دسته از برنامه ها از ابزارهای بتدریج کم رنگ شده EIS و DSS تجاری خاص نزدیکتر و طبیعتا گرانتر هم هستند. البته اختلاف بین خود را (Enterprise) و در سطح سازمان **داده** بتازگی دامنه عمل خود را گسترش EIS است. ابزارهای استفاده می کنند مطرح کرده اند بطوریکه مدیران و تحلیلگران نیز از این ابزارها

(MDA) ابزارهای یا OnLine Analytical Processing (OLAP) بطور خلاصه ابزارهای DSS/ESS ابزارهای. (Data Mining) **کاوی داده** آنها ابزارهای نامیده می شوند و در لایه های بالاتر به Multidimensional Analysis گفته می شود (Intelligent Agent) کارآگاهان شخصی و (Mining)

OLTP و DSS مقایسه سیستم های

را می بینیم OLTP و DSS در جدول زیر تفاوت های دو نوع سیستم

DSS نیاز بانک اطلاعاتی OLTP قابلیت نیاز بانک اطلاعاتی می کند کارکنان سیستم تولید کننده اطلاعات شکارچی اطلاعات چه کسی از آن استفاده مقدار فعلی اطلاعات نیاز دارد و گزارش ها قابل باز سازی نیستند به اطلاعات ارزش زمانی اطلاعات به هستند نیاز دارد . اطلاعات هر از گاه به وقت می شوند. گزارش ها قابل بازسازی پایدار از گاه تعداد دسترسی ها به اطلاعات پیوسته در طول روز کاری با نقاط پیک کاری هر تبدیلی صورت نگرفته در چندین لایه تبدیل صورت گرفته است. خام است. استخراج و **داده داده** شکل ها انجام شده **داده** استخراج و فشرده سازی از چندین محل داخلی و خارجی ها از یک برنامه **داده** آوری جمع شود خیر از برنامه های توسط یک برنامه تولید می **داده** مشخص است بلی بیشتر **داده** آیا محل تولید مختلف و بانک های اطلاعات و وب می آید ها پیوسته و در یک نگارش هستند بلی هر مجموعه **داده** . بندی شده هستند خیر آیا اطلاعات نگارش است دارای تاریخ برداشت **داده** از چندین کاربر اطلاعات را به وقت می کنند بیشتر اوقات یک کاربر **داده** نوع دسترسی به مقدار کنونی مدام در حال تغییر است فقط خواندن نیست قابل به وقت رسانی است **داده** آیا ندارد. فقط از طریق برنامه ها ممکن است. منعطف از طریق یک تولید انعطاف در دسترسی انعطاف OLAP کننده درخواست و نسبتا کند راندمان سرعت پاسخ بالا مورد نیاز است. فعالیت ها همگی مکانیزه و سریع زیادی کار کشف و تحقیق و نیازهای اطلاعاتی بخوبی فهمیده شده اند ناپایدار و نسبی. به مقدار جستجوی موضوعی نیاز است ها ممکن است از هر جایی بیابند **داده** که در بانک موجود است دامنه اطلاعات محدود. آن چیزی هزاران و میلیونها رکورد / رکورد های پردازش شده کمتر از ۱۰ رکورد صدها

داده انبار (Datawarehouse)

اطلاعات برای مصرف سیستم های (Repository) یعنی انبار **داده** در محیط خادم/مخدوم انبار از اطلاعات است که قادر است اطلاعات بک مخزن فعال و هوشمند **داده** پشتیبانی تصمیم گیری. انبار

و نهایتاً پخش نماید و در صورت لزوم نیز سیاست را از محیط های گوناگون جمع آوری و مدیریت کرده نماید های تجاری را روی آنها اجرا

عناصر انبار داری

انباره یک محل است و انباره داری یک فرآیند

: این فرآیند از عناصر زیر تشکیل شده است

روی بانک های مختلف مدیریت انتشار اطلاعات انباره که وظیفه نسخه برداری و توزیع اطلاعات را بر ۱۰ اطلاعاتی را که بایستی کپی (آنگونه که شکارچی اطلاعات تعریف می کند) به عهده دارد. شکارچی تبدیلات لازم روی اطلاعات را تعریف می کند. شود، مبدا و مقصد اطلاعات، تعداد بوقت رسانی ها و کامل آخرین وضعیت اطلاعات و اصطلاح بوقت رسانی بمفهوم کپی (Refresh) اصطلاح تازه سازی بکار گرفته شده اند. همه کارها می تواند بصورت خودکار و یا بمفهوم اعمال آخرین تغییرات (Update) اطلاعات ممکن است از بانکهای رابطه ای و غیر رابطه ای تهیه شود. توجه کنید که دستی انجام پذیرد اطلاعات خارجی قبل از ورود به سیستم، تبدیل شده و پاک سازی می شوند کلیه

نمودن یک بانک اطلاع رسانی یک بانک اطلاعاتی رابطه ایست که وظیفه سازماندهی و ذخیره ۲۰ حاصله از منابع مختلف و نسخه از اطلاعات و همچنین تبدیلات و جمع بندی و افزودن ارزش به اطلاعات در مورد اطلاعات) نیز به عهده این بانک با فرمت های مورد نظر بعهده دارد. نگهداری فراداده (اطلاعات ایندکس ها و غیره را بیان می کنند و فراداده های است. فراداده های سیستمی روابط بین جداول و را برای یک شکارچی اطلاعات روشن می سازند ارزش اطلاعات (semantic) محتوایی

راهنمای تجاری و یک ترکیبی از یک راهنمای فنی و (Informational Directory) راهنمای اطلاعات ۳۰ دانستن محل وجود اطلاعات پوششگر اطلاعات است. هدف اصلی این راهنما کمک به شکارچی برای شکل آن و روش دسترسی به آن است

از فروشندگان انجام می گیرد. بسیاری SQL از طریق انواع دستورات DSS/EIS پشتیبانی ابزارهای ۴۰ و سایرین انواع دیگر پروتکل ها را سرویس می دهند ODBC پروتکل

(DataMarts) داده سلسله مراتب انباره ها (غرفه های

داده دپارتمانی و غرفه های **داده** هستند. در عمل غرفه های **داده** انواع کوچکتری از انباره های ابتدا برنامه ریزی نمی شوند بلکه ابتدایا بوجود آمده و در صورت موفقیت تکثیر شده از (mobile) همراه ها تشکیل در نهایت مدیر بانک اطلاعاتی سازمان ممکن است بتواند یک فدراسیون آزاد از این غرفه و نماید را پایه گذاری **داده** دهد و نهایتاً یک انباره

(DataMining) **کاوی داده** و تا (OLAP) تا تحلیل زنده (Queries) از خواسته ها DSS/EIS ابزارهای

ابزارهای گزارش گیری

خواسته پردازها بما اجازه ساختن یک دستور سیکویل را می دهند بدون آنکه **داده** ابزارهای تحلیل برنامه ای بنویسیم یا سیکویل یاد بگیریم. با چند نشانه و کلیک عبارت های سیکویل مجبور باشیم شود. مناسب برای گرد آوری اطلاعات و نمایش آن بشکل یک گراف / جدول و یا گزارش آماده می را می دهند و به این ابزارهای برجسته تر در این زمینه امکان کنترل میزان نتایج برگشته از یک خواسته را برگردانند گرفت. در سال ۱۹۹۸ ترتیب می توان جلوی درخواستهایی را که ممکن است میلیونها رکورد را از آن جمله اند Microsoft Access, Oracle Reports, Business Objects است که بیش از ۱۵۰ نوع از این ابزارها در بازار وجود داشته

و اطلاعات چند بعدی OLAP

که می توانید آنرا در جهات مختلف **نگاه کنید** ها **داده** مثل یک مکعب روبیک از OLAP به ساختار کنید بچرخانید تا بتوانید سناریو های "قبلا چه شده" و "چه می شد اگر ..." را بررسی

اطلاعاتی دو بعدی (و یا بانکهای خاص) ها را توسط بانکهای **داده** این ابزارها دیدگاههای چند بعدی از قدرت OLAP ها در **داده** توان دسترسی چند بعدی به چند بعدی) تولید کرده و در اختیار ما می گذارند. کردن خواسته های پیچیده تر را بما می دهد فرموله

متعارف فقط یک صفحه گسترده با چند محور است (در صفحات گسترده OLAP برای سادگی فرض کنید داریم) در این صورت مثلا می ... و عمودی با ایندکس های A, B, C, ... دو محور افقی با اختصار فروش، تاریخ، مشتری، فروشگاه، قیمت و توانیم اطلاعات فروش یک سازمان را از دیدگاه های منطقه میزان فروش به ازای یک محصول و فروشگاه در یک میزان فروش بررسی کنیم. و پاسخ سولاتی نظیر ماه مشخص را خواهیم داشت

های اطلاعاتی رابطه ای ها را در مقایسه با بانک **داده** طریقه نمایش دادن OLAP مدل چند بعدی اطلاعاتی رابطه ای سرویس فوق را ارایه با ایجاد یک لایه محافظ روی یک بانک ROLAP. تسهیل می کند ذخیره سازی و محاسبه اطلاعات چند بعدی برای فقط راهی برای OLAP میدهد. از دیدگاه فنی پیش روی چندین محور جمع می ها را از **داده**، OLAP خادم پاسخگوایی به سناریوهای کاربر است. یک **داده** OLAP بایستی پاک سازی شوند. غالبا OLAP زند. توجه کنید که اطلاعات قبل از وارد شدن به استخراج می کند **داده** را از یک انباره

را به چند دسته تقسیم می کنند OLAP ابزارهای

رو میزی OLAP

سازند ابزارهای ساده و مستقل که روی کامپیوتر های شخصی نصب شده و مکعب های کوچکی می صفحات گسترده ای و آنها را نیز بر روی سیستم به شکل فایل ذخیره می کنند. بیشتر این ابزارها با به استفاده از این دسته از کار می کنند. به این ترتیب کسانی که در سفر هستند قادر Excel نظیر کردن این محصولات است در حال جایگزین Web OLAP محصولات هستند. (در حال حاضر

چند بعدی MOLAP

اطلاعاتی خاصی را بجای ذخیره کردن اطلاعات در رکورد های کلید دار، این دسته از ابزارهای بانکهای **داده** مرتب شده بر اساس ابعاد ها را به شکل آرایه های **داده** برای خود طراحی کرده اند بطوریکه در حال حاضر نیز دو استاندارد برای این تیب ابزار وجود دارد. سرعت این (HyperCubes) ذخیره می کنند ابزار بالا ولی سایز بانک اطلاعاتی آن نسبتا کوچک است

رابطه ای OLAP (ROLAP)

می کنند. بطوریکه این ابزار ها با ایجاد یک بستر روی بانکهای رابطه ای اطلاعات را ذخیره و بازیابی بر همین Red Brick, MicroStrategy اساس بینه سازی برخی بانکهای اطلاعاتی رابطه ای مانند اساس استوار است. توجه می باشد اندازه بانک اطلاعاتی این ابزار قابل

Hybrid OLAP (HOLAP)

است (MOLAP طرح شده در) MDBMS و ROLAP ترکیبی از hybrid در اینجا منظور از می باشد ROLAP ابزار دارای بانک اطلاعاتی بزرگ و رادمان بالاتر نسبت به

OLAP استانداردهای

و از طرف دیگر MD-API با استاندارد OLAP با دو استاندارد مواجه است، از یک طرف گروه OLAP جامعه و دومی از حمایت Oracle اولی از حمایت (OLE DB for OLAP (Tensor) با استاندارد Microsoft برنامه های آنان را در ابعاد فروش MS-SQLV فروشندگان کوچکتری برخوردار است که امیدوارند فروش ویندوز مطرح کند

(Data Mining) کاوی داده

طلای کوچکی را که در های ما می توانند تکه **داده** با جستجوی حجم عظیم **کاوی داده** ابزارهای گوشه ای پنهان شده بیابند.

های یک واحد **داده** بازگشت هزینه صرف شده در این ابزارها غالباً بسیار سریع است. مثلاً در بررسی فروشگاه متوجه شدند که میزان سرقت حین فروش از باتریها و فیلمها و قلم های با قیمت ار یک جابجا کردن متوسط ماهانه حدود ۶۰۰۰۰ دلار برای فروشگاه هزینه داشته است که به این ترتیب با جویی بدنبال داشته است اقلام و قرار دادن در قسمتهای با دید بهتر سالانه حدود ۷۰۰۰۰۰ دلار صرفه

گردد که ممکن است از دید ما ها می **داده** بندی هایی در بدنبال طرحها و گروه **کاوی داده** ابزارهای که استفاده OLAP گیرد. بر خلاف ابزارهای پنهان مانده باشد. ابزار تقریباً از کاربر هیچ کمکی نمی استفاده کننده را این ابزار است که **کاوی داده** هستند در کنندگان راهنما و سازمان دهنده اطلاعات که چه می خواهید. بیشتر این هدایت می کند. ابزار فرض می کند که شما خود نیز دقیقاً نمی دانید: ابزارها از روش های جستجوی زیر استفاده میکنند است که وجود ارتباطات که اصطلاحاً تحلیل سید بازار خوانده می شود. ابزار بدنبال اثبات این موضوع ۱- تعطیلات تابستانی در استرالیا چیزی بمعنی وجود چیز دیگریست. مثلاً بیشتر خریداران لوازم غواصی به خریدار کا لای دیگری نیز هست می روند. یا مصرف کننده یک کالای مشخص مصرف کرده

قیمت طلا ۱۰ درصد ارتباطات متوالی ابزار بدنبال روابط متوالی بین موضوعات می گردد مثلاً وقتی ۲- بالا می رود یک هفته بعد قیمت سهام ۱۵ درصد پایین می آید

درصد رای دسته بندی بدنبال دسته بندی و طبقه بندی سطح بالای اطلاعات هستند. مثلاً ۷۰-۳ دارند بین ۴۰ تا ۵۰ سال دهندگانی که تصمیم نگرفته اند به که رای دهند درآمدی بالای ۶۰۰۰۰ دلار اقامت دارند X سن دارند و در منطقه

می رسید افراد اگر اطلاعات جدول زیر در یک گراف دو بعدی به تصویر در آید متوجه می شویم که بنظر بین ۲۳ تا ۲۹ به مکزیک و بین ۳۰ تا ۵۱ به کانادا سفر می کنند

سن مشتری کشوری که به آن سفر کرده

مکزیک ۲۳
کانادا ۴۵
کانادا ۳۲
کانادا ۴۷
کانادا ۴۶
کانادا ۳۴
کانادا ۵۱
مکزیک ۲۸
کانادا ۴۹
مکزیک ۲۹
مکزیک ۲۶
کانادا ۳۱

نمی کنند یک نکته جالب دیگر که بسادگی قابل دیدن نیست آنستکه افراد بین ۳۵ تا ۴۴ اصلاً سفر آنها که بین ۴۵ تا ۵۱ سال عبارت دیگر دو دسته آدم به کانادا سفر می کنند آنها که بین ۳۰ تا ۳۴ و بعدی بسادگی قابل رویت است. چنانچه سن دارند. گروه بندی در این مجموعه اطلاعات کوچک و دو سادگی گذشته نخواهد بود. گفتنی است تعداد نمونه ابعاد اطلاعات و حجم آن افزایش باید موضوع به **کاوی داده** مقادیری که هر یک از ستونها می گیرند در سرعت پردازش ها، تعداد ستونهای اطلاعاتی و پردازش ۱۰۰۰۰۰۰ نمونه با ۲۰۰ ستون اطلاعاتی که هر یک می توانند ۲۵ موثر هستند. مثلاً برای خود بگیرند به حدود ۲ ساعت وقت نیازاست مقدار مختلف به

کاربرد یافته اند. از جمله محققین بهداشت برای کشف میزان موفقیت این ابزارها در زمینه های مختلف برای ارزیابی اعتبار مشتریان، بورس بازان برای تشخیص جابجایی قیمت های سهام و جراحیها، بانکها طرحهای تجاری، شرکتهای بیمه برای تشخیص ریسک مشتریان و رفتارهایشان و هتل ها تشخیص **داده** تشخیص مشتریان بازگشتی خود از آن استفاده میکنند. همانطوریکه بنظر می آید ابزارهای برای رده بالاتر هستند که استفاده های قابل توجهی برای آنها در صنعت قابل از مجموعه ابزارهای یک **کاوی** تصور است.

: برخی از انواع تجاری این ابزار عبارتند از
Intelligent Miner, Darwin, MineSet, KnowledgeStudio, DataMind, Clementine

کارآگاهان شخصی

اطلاعات مورد نیاز را **داده** این مامورین برنامه های قابل حملی هستند که با اتصال به انباره های قوانین تعریف استخراج کرده و به کارفرمایان خود اطلاع می دهند. در حال حاضر این ابزارها بر اساس مشاهده تغییر پیام شده از طرف کارفرمای خود به جستجوی تغییرات در اطلاعات رفته و در صورت مناسب را می دهند

: هنوز کار های زیادی در این قسمت بایستی صورت پذیرد که از آن جمله اند ابزار (هوشمند شدن)، درک علایق کارفرما و جستجو دربانکهای ها بر اساس دانش درون **داده** درک اعلام تغییرات به کارفرما اطلاعاتی مختلف برای کاربر و یا برنامه های همسر یابی با توجه به برنامه های کاریابی روی اینترنت با توجه به رزومه و سایر شرکتها اعلام می کنند و Microsoft سایتهاى مشخصات. برنامه هایی که تغییرات را در هستند از این نوع برنامه ها **مثالهای ساده ای**

[کپی رایت و منابع](#)

: مقالات مرتبط

داشته باشد، موجود نیست **مقاله** ارتباط زیادی با این ای که **مقاله** در حال حاضر متأسفانه

منتشر **کامپیوتر** ب.ظ و درباره موضوعات ۳:۴۶ این مطلب در تاریخ سه شنبه ۹ آبان ۱۳۸۵ در ساعت **نظر** میتوانید دنبال کنید. شما **RSS 2.0** شما میتوانید هر پاسخی به این مطلب را توسط .شده است **تالارهای گفتگو** لطفاً سوالات و درخواست های خود را در .ارسال کنید **دنبالک** ، یا از سایت خود **بدهید** .نمائید مطرح

شده است داده " (Datawarehouse) داده یک نظر درباره "آشنایی با مفاهیم انباره های

: گفت **محسن بختیاری** 1.
[ساعت ۱۲:۵۹ ب.ظ چهارشنبه ۲۶ اردیبهشت ۱۳۸۶](#)

تشکر ای می خواستم .با **مقاله** در مورد فایلهاى سیستمی

معرفی داده کاوی

داده کاوی [1] یا کشف دانش در پایگاه داده ها [2] (**KDD**) علم نسبتاً تازه ای است که با توجه پیشرفت کشور در زمینه IT و نگاه های ویژه به دولت الکترونیک و نفوذ استفاده از سیستم های رایانه ای در صنعت و ایجاد بانک های اطلاعاتی بزرگ توسط ادارات دولتی، بانک ها و بخش خصوصی نیاز به استفاده از آن به طور عمیقی احساس می شود. داده کاوی یعنی کشف دانش و اطلاعات معتبر پنهان در پایگاه های داده. یا به بیان بهتر تجزیه و تحلیل ماشینی داده ها برای پیدا کردن الگوهای مفید و تازه و قابل استناد در پایگاه داده های بزرگ ، داده کاوی نامیده می شود. داده کاوی در پایگاه های داده کوچک نیز بسیار پرکاربرد است و از نتایج و الگوهای تولید شده بوسیله آن در تصمیم گیری های استراتژیک تجاری شرکت های کوچک نیز می توان بهره های فراوان برد. کاربرد داده کاوی در یک جمله را این گونه می توان بیان کرد : " داده کاوی اطلاعاتی می دهد ، که شما برای گرفتن تصمیم هوشمندانه ای درباره مشکلات سخت شغلان به آنها نیاز دارید" [3] .

مثالی کلاسیک از کاربرد داده کاوی

اغلب تجارت ها به تصمیم گیری های استراتژیک و یا اتخاذ خط مشی های جدید برای خدمت رسانی بهتر به مشتریان نیاز دارند. به عنوان مثال فروشگاهها آرایش مغازه خود را برای ایجاد میل بیشتر به خرید مجدداً طراحی می کنند. این مثال به داده هایی در مورد رفتار مصرفی گذشته مشتریان برای تعیین الگوهای به وسیله داده کاوی، نیاز دارند.

برای روشن تر شدن مسئله می توان مثال را اینگونه بیان کرد که در یک فروشگاه زنجیره ای پس از داده کاوی مشخص میشود که درصدی از مشتریان خرید تلویزیون ، میز تلویزیون و گلدان کریستالی را هم در همان روز و بعد از خرید تلویزیون میخرند.مدیر فروشگاه می تواند بلافاصله دستوراتی صادر کند که براساس مدل های تلویزیون موجود میزهایی و براساس مدل میزها گلدانهای کریستالی برای فروش سفارش داده شود و غرفه های جنبی غرفه تلویزیون را به میز و گلدان کریستالی اختصاص دهد. مطمئناً حتی پس از مدت کوتاهی سود حاصل از این بخش از فروشگاه به طور قابل ملاحظه ای ترقی خواهد کرد.

در واقع ابزار داده کاوی، داده را می گیرد و يك تصویر از واقعیت به شکل مدل می سازد، این مدل روابط موجود در داده ها را شرح می دهد.

برای بهبود بهره وری از یک فروشگاه داده کاوی از داده های انبار داده ، مدل هایی را ارائه میدهد که بیانگر این هستند که چه محصولات یا خدماتی، به چه مشتریانی، در چه زمانی و از طریق چه کانالی عرضه شود.

بیشتر شرکتها، بانکهای داده ای عظیمی شامل داده های بازاریابی، منابع انسانی و مالی را دارا هستند. بنابراین، سرمایه گذاری در زمینه انبار داده، یکی از اجزای حیاتی در استراتژی مدیریت ارتباط با مشتری است.

رابطه مشتری با زمان تغییر می کند و چنانچه تجارت و مشتری درباره یکدیگر بیشتر بدانند این رابطه تکامل و رشد می یابد. چرخه زندگی مشتری چارچوب خوبی برای به کارگیری داده کاوی در مدیریت ارتباط با مشتری فراهم می کند. در بخش ورودی داده کاوی، چرخه زندگی مشتری می گوید چه اطلاعاتی در دسترس است و در بخش خروجی آن، چرخه زندگی می گوید چه چیزی احتمالاً جالب توجه است و چه تصمیماتی باید گرفته شود. داده کاوی می تواند سودآوری مشتری های بالقوه را که می توانند به مشتریان بالفعل تبدیل شوند، پیش بینی کند و اینکه تا چه مدت به صورت مشتریان وفادار خواهند ماند و چگونه احتمالاً ما را ترک خواهند کرد.

بعضی از مشتریان مرتباً مراجعاتشان را به شرکتها برای کسب مزیتهایی که طی رقابت میان آنها به وجود می آید، تغییر می دهند. در این صورت شرکتها می توانند هدفشان را روی مشتریانی متمرکز کنند که سودآوری بیشتری دارند.

بنابراین می توان از طریق داده کاوی ارزش مشتریان را تعیین، رفتار آینده آنها را پیش بینی و تصمیمات آگاهانه ای را در این رابطه اتخاذ کرد.

از کاربرد های داده کاوی می توان به نمونه های زیر اشاره کرد :

۱. بانکداری :

- از جالب توجه ترین کاربردهای داده کاوی می توان به کشف پول شویی اشاره کرد.
- تشخیص مشتریان ثابت و همیشگی
- تعیین مشتریان استفاده کننده از یک سرویس خاص

۲. بیمه :

- پیش گویی میزان استقبال از بیمه نامه های جدید
- تشخیص کلاهبرداری ها و مشخص کردن رفتار های نامتناسب
- تشخیص نیاز مشتریان و خواسته های آنها
- تشخیص تخلفات پزشکی

واضح است که زمینه استفاده از داده کاوی بی نهایت گسترده است، و دو مثال فوق به خاطر درک راحت تر انتخاب شده اند.

داده کاوی شباهت زیادی به تحلیل های آماری دارد، ولی داده کاوی از جهات زیادی با آمار متفاوت است و مزیت های زیادی نسبت به آمار دارد. جالب ترین تفاوت داده کاوی با تحلیل های آماری این است که در آمار ما فرضیه ای طرح می کنیم و با استفاده از تحلیل های آماری به اثبات یا رد فرضیه می پردازیم اما داده کاوی به فرضیه احتیاجی ندارد. در واقع ابزار داده کاوی فرض می کند که شما خود هم نمی دانید به دنبال چه می گردید. و این نکته ای است که باعث می شود کار آمدی داده کاوی در مواقع بروز مشکل نمایان شود. برای مثال ما در آمار فرض می کنیم که دو گروه فاصله ای باهم ارتباط دارند سپس با استفاده از ضریب هم بستگی پیرسون مشخص می کنیم که ارتباط وجود دارد یا خیر. ولی داده کاوی بدون توجه به اینکه ما اینگونه فرضی داشته باشیم یا نه با کاوش میان داده ها اگر ارتباطی مخفی معنی داری وجود داشته باشد آن را به اطلاع ما می رساند. تفاوت بعدی آمار و داده کاوی در این است که آمار فقط می تواند از داده های عددی استفاده کند ولی داده کاوی از داده های غیر عددی هم استفاده می کند. تفاوت های دیگری هم میان آمار و داده کاوی وجود دارد که بحث در مورد آنها در حوصله این مقاله نمی گنجد.

اما برای اولین بار در سال ۱۹۵۰ از رایانه برای تحلیل و ذخیره پایگاه داده‌ها استفاده شد. ولی حجم اطلاعات و میزان رشد آنها به قدری زیاد بوده است که هم اکنون کسی از میزان اطلاعات ذخیره شده در پایگاه داده‌های سراسر دنیا به صورت دقیق اطلاعی ندارد ولی مطمئناً حجم اطلاعات و مخصوصاً سرعت رشد آنها به قدری زیاد شده که آمار شناسان و تحلیل‌گران در بررسی و تحلیل پایگاه‌های داده در زمینه‌های مختلف ناتوانند. بعضی از پایگاه داده‌ها به قدری بزرگ و پیچیده شده‌اند که تحلیل روابط و استخراج اطلاعات مفید پنهان شده در آنها واقعا از ظرفیت ذهنی بشری فراتر رفته است. از زمانی که رشد پایگاه‌های داده و حجم اطلاعات، سرعت گرفت و میزان داده‌ها افزایش یافت، نیاز به تحلیل ماشینی داده‌ها و استخراج سریع و دقیق دانش نهفته در آنها احساس شد. شاید بتوان لوول (۱۹۸۳) را اولین شخصی دانست که گزارشی در مورد داده کاوی تحت عنوان «شبهه سازی فعالیت داده کاوی» ارائه نمود. [4]

عمل داده کاوی از یک پایگاه داده به چند مرحله مشخص تقسیم می‌شود که ما در این مقاله به معرفی و توضیح مختصر در مورد هر یک از این مراحل اکتفا می‌کنیم:

۱. مرحله اول: تشکیل انبار داده .
با توجه به عنوان، این مرحله برای تشکیل محیطی پیوسته و یک پارچه جهت انجام مراحل بعدی و داده کاوی در آن، انجام می‌گیرد. در حالت کلی انبار داده مجموعه پیوسته و طبقه بندی شده است که دائماً در حال تغییر بوده و دینامیک است که برای کاوش آماده می‌شود.
۲. مرحله دوم: انتخاب داده‌ها
در این مرحله برای کم کردن هزینه‌های عملیات داده کاوی، داده‌هایی از پایگاه داده انتخاب می‌شوند که مورد مطالعه هستند و هدف داده کاوی دادن نتایجی در مورد آنهاست.
۳. مرحله سوم: تبدیل داده‌ها .
مشخص است برای انجام عملیات داده کاوی لزوماً باید تبدیلات خاصی روی داده‌ها انجام گیرد ممکن است این تبدیلات خیلی راحت و مختصر مثل تبدیل byte به integer باشد یا خیلی پیچیده و زمان‌بر و با هزینه‌های بالا مثل تعریف صفات جدید و یا تبدیل و استخراج داده‌ها از مقادیر رشته‌ای و ... باشد.
۴. مرحله چهارم: کاوش در داده‌ها .
در این مرحله است که داده کاوی انجام می‌شود. در این مرحله با استفاده از تکنیک‌های داده کاوی داده‌ها مورد کاوش قرار گرفته، دانش نهفته در آنها استخراج شده و الگو سازی صورت می‌گیرد.
۵. مرحله پنجم: تفسیر نتیجه .
در این مرحله نتایج و الگوهای ارائه شده توسط ابزار داده کاو مورد بررسی قرار گرفته و نتایج مفید معین می‌شود.

طرز کار ابزار داده کاو اینگونه است که ابزار به دنبال اثبات این است که وجود چیزی به معنای وجود چیز دیگری است و سعی می‌کند در درجه اول از توالی ارتباطات برای کشف یک الگو بهره بگیرد و در نهایت اطلاعات بدست آمده را دسته بندی کند تا به الگوی خاصی برسد که بتواند آن را براساس فاکتورهای داخلی به مخاطبش ارائه دهد.

همچنین در داده کاوی از الگوریتم‌های ژنتیک و شبکه‌های عصبی هم استفاده می‌شود. شبکه‌های عصبی به علت کار آمدی در حل مسائل پیچیده و بزرگ مورد استفاده‌اند و کاربرد الگوریتم‌های ژنتیک

در داده کاوی برای جستجو و ساختن یک مدل بهینه در میان مدل های بدست آمده است ، به این گونه که مدل های اولیه روی کروموزوم هایی قرار می گیرند و با رقابت بر سر انتقال صفات به نسل بعد ، بهترین مدل و لایق ترین آنها به کاربر ارائه می شوند.

داده کاوی امروز گسترش زیادی یافته است به طوری که اکثر نرم افزار های پایگاه داده ای مثل SQL Server و ORACLE نیز شامل ابزارهایی داده کاوی شده اند ولی هنوز نرم افزار های تخصصی داده کاوی همچون Intelligent Miner , Darwin , Mine Set, Knowledge Studio, Data Mind از مهمترین ابزار های داده کاوی اند.

منابع

1-CHRIS RYGIELSKI, "DATA MINING TECHNIQUES FOR CUSTOMER RELATIONSHIP MANAGEMENT", TECHNOLOGY IN SOCIETY, 2002 .

2- HILL L., "CRM: EASIER SAID THAN DONE", INTELLIGENT ENTERPRISE, 1999

4- Microsoft Visual Studio .Net Documentation

5- Client/Server Survival Guide by Robert Orfali, Dan Harkey, Jeri Edwards

۶- شاه سمندی، پرستو «داده کاوی در مدیریت ارتباط با مشتری» (۱۳۸۴)، مجله تدبیر شماره ۱۵۶.

7- Hand. D.J (1998): "Review of Data mining", The American statistician, 52, 112-118.

8- Jeffery W. Seifert , Analyst in information science and Technology Policy, ` Data Mining : An Overview ` December 2004.

داده کاوی پل ارتباطی میان علم آمار، علم کامپیوتر، هوش مصنوعی، الگوشناسی، فراگیری ماشین و بازنمایی بصری داده می باشد. داده کاوی فرآیندی پیچیده جهت شناسایی الگوها و مدل های صحیح، جدید و به صورت بالقوه مفید، در حجم وسیعی از داده می باشد، به طریقی که این الگوها و مدلها برای انسانها قابل درک باشند. داده کاوی به صورت یک محصول قابل خریداری نمی باشد، بلکه یک رشته علمی و فرآیندی است که بایستی به صورت یک پروژه پیاده سازی شود.



داده ها اغلب حجیم می باشند و به تنهایی قابل استفاده نیستند، بلکه دانش نهفته در داده ها قابل استفاده می باشد. بنابراین بهره گیری از قدرت فرآیند داده کاوی جهت شناسایی الگوها و مدلها و نیز ارتباط عناصر مختلف در پایگاه داده جهت کشف دانش نهفته در داده ها و نهایتا تبدیل داده به اطلاعات، روز به روز ضروری تر می شود.

★ خدمات شرکت دانا پرداز در این زمینه

مثال تفهیمی در مورد داده کاوی

یکی از نمونه های بارز داده کاوی را می توان در فروشگاه های زنجیره ای مشاهده نمود، که در آن سعی می شود ارتباط محصولات مختلف هنگام خرید مشتریان مشخص گردد. فروشگاه های زنجیره ای مشتاقند بدانند که چه محصولاتی با یکدیگر به فروش می روند. برای مثال طی یک عملیات داده کاوی گسترده در یک فروشگاه زنجیره ای در آمریکا شمالی که بر روی حجم عظیمی از داده های فروش صورت گرفت، مشخص گردید که مردانی که برای خرید فنجان بچه به فروشگاه می روند معمولا آب جو نیز خریداری می کنند. همچنین مشخص گردید مشتریانی که تلویزیون خریداری می کنند، غالبا گلدان کریستالی نیز می خرند. نمونه مشابه عملیات داده کاوی را می توان در یک شرکت بزرگ تولید و عرضه پوشاک در اروپا مشاهده نمود، به شکلی که نتایج داده کاوی مشخص می کرد که افرادی که کراوات های ابریشمی خریداری می کنند، در همان روز یا روزهای آینده گیره کراوات مشکی رنگ نیز خریداری می کنند. به روشنی این مطلب قابل درک است که این نوع استفاده از داده کاوی می تواند فروشگاه ها را در برگزاری هوشمندانه فستیوال های فروش و نحوه ارائه اجناس به مشتریان یاری رساند. نمونه دیگر استفاده از داده کاوی در زمینه فروش را می توان در یک شرکت بزرگ دوبلاژ و تکثیر و عرضه فیلم های سینمایی در آمریکا شمالی مشاهده نمود که در آن عملیات داده کاوی، روابط مشتریان و هنرپیشه های سینمایی و نیز گروه های مختلف مشتریان بر اساس سبک فیلم ها (ترسناک، رمانتیک، حادثه ای و ...) مشخص گردید. بنابراین آن شرکت به صورت کاملا هوشمندانه می توانست مشتریان بالقوه فیلم های سینمایی را بر اساس علاقه مشتریان به هنرپیشه های مختلف و سبک های سینمایی شناسایی کند.

از دیگر زمینه های به کارگیری داده کاوی، استفاده بیمارستانها و کارخانه های داروسازی جهت کشف الگوها و مدلهای ناشناخته تاثیر دارو ها بر بیماری های مختلف و نیز بیماران گروه های سنی مختلف را می توان نام برد. استفاده از داده کاوی در زمینه های مالی و بانکداری به شناخت مشتریان پر خطر و سودجو بر اساس معیارهایی از جمله سن، درآمد، وضعیت سکونت، تحصیلات، شغل و غیره می انجامد.

تعاریف داده کاوی

داده کاوی استخراج اطلاعات مفهومی، ناشناخته و به صورت بالقوه مفید از پایگاه داده می باشد.

Source: W.Frawley and G. Piatetsky. Knowledge Discovery I DataBases.ISSN 0738-4602

داده کاوی علم استخراج اطلاعات مفید از پایگاه های داده یا مجموعه داده ای می باشد.
Source: D. Hand,H. Mannila,P. Smyth(2001).Principles of Data Mining.MIT Press,Cambridge

داده کاوی استخراج نیمه اتوماتیک الگوها، تغییرات، وابستگی ها، نابهنجاری ها و دیگر ساختارهای معنی دار آماری از پایگاه های بزرگ داده می باشد .

Source: R.Grossman

خدمات شرکت دانا پرداز در این زمینه ★

تفاوت داده کاوی و آنالیز های آماری



داده کاوی معمولا با نوشتن مقدار زیادی گزارش و تحقیق و استعلام در آنها اشتباه گرفته می شود. اما در واقع داده کاوی هیچ کدام از اینها را شامل نمی شود. داده کاوی توسط تجهیزات خاصی صورت می پذیرد، که عملیات کاوش را بر اساس تجزیه و تحلیل مکرر داده ها انجام می دهد .

داده کاوی با آنالیز های متداول آماری نیز متفاوت است؛ در زبیری توان برخی از اصلی ترین تفاوت های داده کاوی و آنالیز آماری را مشاهده نمود :

آنالیز آماری:

- آمار شناسان همیشه با یک فرضیه شروع به کار می کنند.
- آنها از داده های عددی استفاده می کنند.
- آمارشناسان باید رابطه هایی را ایجاد کنند که به فرضیه آنها مربوط است.
- آنها می توانند داده های نابجا و نادرست را در طول آنالیز مشخص کنند.
- آنها می توانند نتایج کار خود را تفسیر و برای مدیران بیان کنند.

کاوی:

داده

- به فرضیه احتیاجی ندارد.
- ابزارهای داده کاوی از انواع مختلف داده ، نه تنها عددی می توانند استفاده کنند.
- الگوریتمهای داده کاوی به طور اتوماتیک روابط را ایجاد می کنند.

• داده کاوی به داده های صحیح و درست نیاز دارد.

• نتایج داده کاوی نسبتاً پیچیده می باشد و نیاز به متخصصانی جهت بیان آنها به مدیران دارد.

جهت درک بهتر تفاوت داده کاوی و آنالیزهای آماری به مثال زیر که در مورد شناخت کلاهبرداری های شرکت بیمه می باشد، توجه کنید.

روش آنالیز آماری:

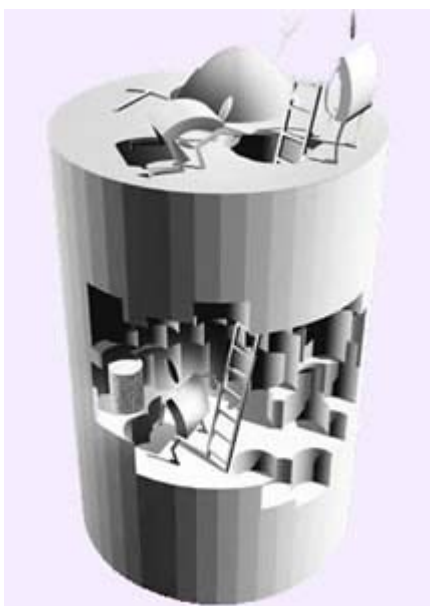
یک مفسر ممکن است متوجه الگوی رفتاری شود که سبب کلاهبرداری بیمه گردد. بر اساس این فرضیه، مفسر به طرح یک سری سوال می پردازد تا این موضوع را بررسی کند. اگر نتایج حاصله مناسب نبود، مفسر فرضیه را اصلاح می کند و یا با انتخاب فرضیه دیگری مجدداً شروع می کند. این روش نه تنها وقت گیر است بلکه به قدرت تجزیه و تحلیل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

روش داده کاوی:

یک مفسر سیستم های داده کاوی را ساخته و پس از طی مراحل از جمله جمع آوری داده ها، یکپارچه سازی و خلاص سازی داده ها به انجام عملیات داده کاوی می پردازد. داده کاوی تمام الگوهای غیرعادی را که از حالت عادی و نرمال انحراف دارند و ممکن است منجر به کلاهبرداری شوند را پیدا می کند. نتایج داده کاوی حالت های 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.

★ خدمات شرکت دانا پرداز در این زمینه

فوائد و نقش داده کاوی در فعالیت شرکتها



امروزه عملیات داده کاوی به صورت گسترده توسط تمامی شرکت هایی که مشتریان در کانون توجه آنها قرار دارند، استفاده می شود، از جمله فروشگاه ها، شرکت های مالی، ارتباطاتی، بازاریابی و غیره .

استفاده از داده کاوی به این شرکتها کمک می کند تا ارتباط عوامل داخلی از جمله قیمت، محل قرارگیری محصولات، مهارت کارمندان را با عوامل خارجی از جمله وضعیت اقتصادی، رقابت در بازار و محل جغرافیایی مشتریان کشف نمایند .

از آنجائیکه هوش مصنوعی یکی از اصلی ترین عناصر داده کاوی

می باشد و با توجه به اینکه به کمک سیستم های کامپیوتری و پایگاه های داده، روزانه به میزان داده ها افزوده می شود، بنابراین استفاده هوشمندانه از دانش بالقوه ای که در این داده نهفته است در دنیای رقابتی امروز برای شرکت ها حیاتی می باشد .

داده کاوی پیش بینی وضع آینده بازار، گرایش مشتریان و شناخت سلیقه های عمومی آنها را برای شرکت ها ممکن می سازد .

مراحل اصلی داده کاوی

داده کاوی را " کشف دانش در داده ها " نیز می نامند. کشف دانش داده ها دارای مراحل مختلفی می باشد که در اینجا به صورت خلاصه آنها را بیان می کنیم :

- استخراج اطلاعات از چندین منبع داده (پایگاه داده).
- یکپارچه سازی اطلاعات و حذف داده های زاید.
- قرار دادن اطلاعات اصلاح شده در انبار داده ها.
- انجام عملیات داده کاوی توسط نرم افزار های مخصوص.
- نمایش نتایج به صورت قابل فهم مانند گزارش و گراف.

انبار داده (Data Warehousing) چیست؟



انبار داده به مجموعه ای از داده

ها گفته

می شود که از منابع مختلف اطلاعاتی سازمان جمع آوری ، دسته بندی و ذخیره می شود. در واقع یک انبار داده مخزن اصلی کلیه داده های حال و گذشته یک سازمان می باشد که برای همیشه جهت انجام عملیات گزارش گیری و آنالیز در دسترس مدیران می باشد . انبارهای داده حاوی داده هایی هستند که به مرور زمان از سیستم های عملیاتی آنلاین سازمان (OLTP) استخراج می شوند، بنابراین سوابق کلیه اطلاعات و یا بخش عظیمی از آنها را می توان در انبار داده ها مشاهده نمود .



از آنجائیکه انجام عملیات آماری و گزارشات پیچیده دارای بارکاری بسیار سنگینی برای سرورهای پایگاه داده می باشند، وجود انبار داده سبب می گردد که اینگونه عملیات تاثیری بر فعالیت برنامه های کاربردی سازمان (OLTP) نداشته باشد.

همانگونه که پایگاه داده سیستمهای عملیاتی سازمان (برنامه های کاربردی) به گونه ای طراحی می شوند که انجام تغییر و حذف و اضافه داده به سرعت صورت پذیرد، در مقابل انبار داده ها دارای معماری ویژه ای می باشند که موجب تسریع انجام عملیات آماری و گزارش گیری می شود (OLAP) .

★ خدمات شرکت دانا پرداز در این زمینه

تاریخچه و دلایل استفاده از انبار داده

از اواخر سال ۱۹۸۰ میلادی، انبار های داده به عنوان نوع متمایزی از پایگاه های داده مورد استفاده اغلب سازمانها و شرکت های متوسط و بزرگ واقع شدند. انبار های داده جهت رفع نیاز رو به رشد مدیریت داده ها و اطلاعات سازمانی که توسط پایگاه های داده سیستم های عملیاتی غیر ممکن بود، ساخته شدند.

سیستمهای عملیاتی سازمان (OLTP) دارای نقاط ضعفی می باشند که انبار های داده آنها را رفع می کنند. از جمله:

- بار پردازش گزارشات موجب کندی عملکرد برنامه های کاربردی می گردد.
- پایگاه های داده برنامه های کاربردی دارای طراحی مناسبی جهت انجام عملیات آماری و گزارش نیستند.
- بسیاری از سازمانها دارای بیش از یک برنامه کاربردی (منابع اطلاعاتی) می باشند، بنابراین تهیه گزارشات در سطح سازمان غیر ممکن می شود.

- تهیه گزارشات در سیستمهای عملیاتی غالباً نیازمند نوشتن برنامه های مخصوص می باشد که معمولاً کند و پرهزینه هستند.

★ خدمات شرکت دانا پرداز در این زمینه

مراحل و نحوه ایجاد انبار داده در سازمان



بسیاری از شرکت ها و سازمانها به این باور رسیده اند که گردآوری، سازمان دهی و یکپارچه سازی داده ها در یک مخزن داده برای مدیریت بهینه و اتخاذ تصمیمات کلان یک ضرورت می باشد .

به طور کلی ساخت یک انبار داده، به شکل یک پروژه شامل مراحل اصلی زیر می باشد:

۱- استخراج داده های تراکنشی از پایگاه های داده به یک مخزن واحد

شناسخت منابع داده های سازمان و استخراج داده های ارزشمند از آنها یکی از اصلی ترین مراحل ایجاد انبار داده می باشد .

۲- تبدیل داده ها

از آنجائیکه سیستمهای اطلاعاتی و برنامه های کاربردی یک سازمان غالباً توسط افراد و پروژه های مختلف به مرور زمان در مواجهه با نیازهای جدید ساخته یا تغییر شکل داده می شوند، یکسان سازی آنها امری ضروری می باشد. در بسیاری از موارد نیز سیستمهای اطلاعاتی در بسترهای مختلف پایگاه داده مانند Microsoft SQL Server، Oracle، Sybase، Microsoft Access و غیره طراحی گردیده اند. بررسی جداول، برقراری ارتباط بین فیلهها و یک شکل سازی داده ها در این مرحله صورت می پذیرد.

۳- بارگذاری داده های تبدیل شده به یک پایگاه داده چند بعدی

بر خلاف پایگاه داده سیستمهای عملیاتی (OLTP) که دارای معماری رابطه ای می باشند و از اصول نرمالیزه استفاده می کنند، طراحی انبار داده به شکلی ویژه بدون بهره گیری از اصول نرمالیزاسیون می باشد. در انبار داده فیلهها در جاهای مختلفی تکرار می شوند و روابط بین جداول کمتر به چشم می خورد. علت آن هم افزایش سرعت پردازش اطلاعات هنگام گزارشات و عملیات آماری می باشد.

۴- تولید مقادیر از پیش محاسبه شده جهت افزایش سرعت گزارش گیری

مقادیر از پیش محاسبه شده را تراکم نیز می نامند. این مرحله توسط سیستمهایی نظیر Microsoft SQL Server Analysis Services بسیار ساده تر شده است. این تراکم ها که در ابعاد مختلف انبار داده ساخته می شوند، موجب می شوند که سرعت انجام عملیات گزارش گیری به شکل محسوسه افزایش یابد. باید توجه داشت که عملیات ساخت این مقادیر بسیار زمان گیر بوده و نیازمند حافظه زیادی بر روی سرور است.

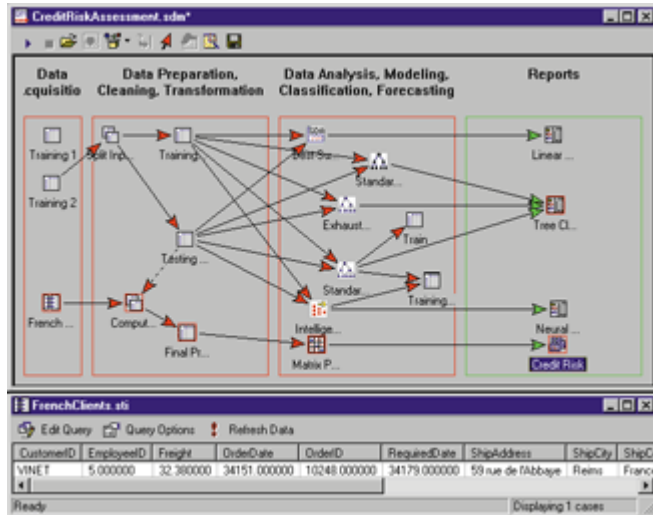
۵- ساخت (یا خرید) یک ابزار گزارش گیری

پس از انجام مراحل فوق، شما می توانید نسبت به ساخت یا خرید یک نرم افزار گزارش گیری تصمیم گیری نمایید. به طور معمول هزینه ساخت یک نرم افزار گزارش گیری، بالاتر از هزینه خرید آن از یک شرکت خارجی می شود.

کاوی داده

داده کاوی فرایندی تحلیلی است که برای کاوش داده ها (معمولاً حجم عظیمی از داده ها - در زمینه های کسب و کار و بازار) صورت می گیرد و یافته ها با بهره گیری از الگوهایی، احراز اعتبار می شوند. هدف اصلی داده کاوی پیش بینی است. فرایند داده کاوی شامل سه مرحله می باشد: ۱. کاوش اولیه، ۲. ساخت مدل یا شناسایی الگو یا کمک احراز اعتبار/ تایید و ۳. بهره برداری.

مرحله ۱: کاوش. معمولاً این مرحله با آماده سازی داده ها صورت می گیرد که ممکن است شامل پاک سازی داده ها، تبدیل داده ها و انتخاب زیرمجموعه هایی از رکوردها با حجم عظیمی از متغیرها (فیلهها) باشد. سپس با توجه به ماهیت مساله تحلیلی، این مرحله به مدل های پیش بینی ساده یا مدل های آماری و گرافیکی برای شناسایی متغیرهای مورد نظر و تعیین پیچیدگی مدل ها برای استفاده در مرحله بعدی نیاز دارد.



مرحله ۲: ساخت و احراز اعتبار مدل. این مرحله به بررسی مدل‌های مختلف و گزینش بهترین مدل با توجه به کارایی پیش‌بینی آن می‌پردازد. شاید این مرحله ساده به نظر برسد، اما اینطور نیست. تکنیک‌های متعددی برای رسیدن به این هدف توسعه یافتند. و "ارزیابی رقابتی مدل‌ها" نام گرفتند. بدین منظور مدل‌های مختلف برای مجموعه داده‌های یکسان به کار می‌روند تا کارایی‌شان با هم مقایسه شود، سپس مدلی که بهترین کارایی را داشته باشد، انتخاب می‌شود. این تکنیک‌ها عبارتند از: Stacking, Bagging, Boosting و Meta-learning.

مرحله ۳: بهره برداری. آخرین مرحله مدلی را که در مرحله قبل انتخاب شده است، در داده‌های جدید به کار می‌گیرند تا پیش‌بینی‌های خروجی‌های مورد انتظار را تولید نماید. داده کاوی به عنوان ابزار مدیریت اطلاعات برای تصمیم‌گیری، عمومیت یافته است. اخیراً، توسعه تکنیک‌های تحلیلی جدید در این زمینه مورد توجه قرار گرفته است (مثلاً *Classification Trees*). اما هنوز داده کاوی مبتنی بر اصول آماری نظیر *Exploratory Data Analysis (EDA)* می‌باشد.

با این وجود تفاوت عمده‌ای بین داده کاوی و EDA وجود دارد. داده کاوی بیشتر به برنامه‌های کاربردی گرایش دارد تا ماهیت اصلی پدیده. به عبارتی داده کاوی کمتر با شناسایی روابط بین متغیرها سروکار دارد.

مفاهیم اساسی در داده کاوی

Bagging: این مفهوم برای ترکیب رده بندی‌های پیش‌بینی شده از چند مدل به کار می‌رود. فرض کنید که قصد دارید مدلی برای رده بندی پیش‌بینی بسازید و مجموعه داده‌های مورد نظرتان کوچک است. شما می‌توانید نمونه‌هایی (با جایگزینی) را از مجموعه داده‌ها انتخاب و برای نمونه‌های حاصل از درخت رده بندی (مثلاً RT&C و CHAID) استفاده نمایید. به طور کلی برای نمونه‌های مختلف به درخت‌های متفاوتی خواهید رسید. سپس برای پیش‌بینی با کمک درخت‌های متفاوت به دست آمده از نمونه‌ها، یک رای‌گیری ساده انجام دهید. رده بندی نهایی، رده بندی‌ای خواهد بود که درخت‌های مختلف آنرا پیش‌بینی کرده‌اند.

Boosting: این مفهوم برای تولید مدل‌های چندگانه (برای پیش‌بینی یا رده بندی) به کار می‌رود. Boosting نیز از روش RT&C یا CHAID استفاده و ترتیبی از classifierها را تولید خواهد کرد.

Meta-Learning: این مفهوم برای ترکیب پیش‌بینی‌های حاصل از چند مدل به کار می‌رود. و هنگامی که انواع مدل‌های موجود در پروژه خیلی متفاوت هستند، کاربرد دارد. فرض کنید که پروژه داده کاوی شما شامل Tree classifierها (مثلاً RT&C و CHAID)، تحلیل خطی و شبکه‌های عصبی است. هر یک از کامپیوترها، رده بندی‌هایی را برای نمونه‌های پیش‌بینی کرده‌اند. تجربه نشان می‌دهد که ترکیب پیش‌بینی‌های چند روش دقیق‌تر از پیش‌بینی‌های هر یک از روش‌هاست. پیش‌بینی‌های حاصل از چند classifier را می‌توان به عنوان ورودی meta-linear مورد استفاده قرار داد. meta-linear پیش‌بینی‌ها را ترکیب می‌کند تا بهترین رده بندی پیش‌بینی شده حاصل شود.

یکپارچه سازی، داده کاوی و طراحی مخازن داده

توسعه بکارگیری سیستم‌های اطلاعاتی در سازمانها و گسترش انتقال فرآیندهای سازمانی به سامانه‌های الکترونیکی و نیز وابستگی روزافزون مدیران و تصمیم‌سازها به اطلاعات و سیستم‌های اطلاعاتی، مباحث نوینی نظیر یکپارچه سازی (Integration)، داده کاوی (Data mining) و مخزن داده (Warehouse Data) را مطرح ساخته است.

مدیران و تصمیم‌گیران سازمانها در وهله اول نیازمند دسترسی به داده‌ها در هر محل و در مرحله بعد به تحلیل آنها می‌باشند تا بتوانند به مزیت رقابتی در بازار در مقابل رقیبان دست یابند. این تحلیل‌ها شامل تشخیص خوشه‌بندی در داده و یافتن تمایلات مشتریان است که از حوزه پایگاه داده‌های معمولی خارج می‌باشد. برای انجام این مهم نه تنها داده بلکه سابقه آن نیز ضروری است. مخزن داده برای اینگونه تحلیل‌ها قادر است داده را از منابع مختلف روی سیستم عامل‌های متفاوت و غیره جمع‌آوری و خلاصه‌بندی کند و بعد از فراهم سازی مخزن داده ابزارهایی برای کاربر نهایی باید فراهم گردد تا بتواند بنحو مطلوب و دلخواه از آن استفاده کند. پردازش تحلیلی برخط (OLAP) و داده کاوی این ابزارها را فراهم می‌کند. در این کارگروه عملیات تحلیل و طراحی مخزن داده و نهایتاً پیاده‌سازی آن انجام گرفته، ابزارهای مختلف داده کاوی و OLAP برای تحلیل تدوین می‌شود. این کارگروه تاکنون تجربه کاری متعدد درباره OLAP با استفاده از امکانات سامانه Oracle z و همچنین تجربه پژوهشی در زمینه داده کاوی عمدتاً بر پایه مدلسازی پیش‌بینی مبتنی بر سیستم‌های یادگیری را دارا می‌باشد.

هم اکنون در هر کشور، سازمان یا شرکت برای امور بازرگانی، پرسنلی، آموزشی، آماری و ... پایگاه داده‌ها ایجاد یا خریداری شده است، به طوری که این پایگاه داده‌ها برای مدیران، برنامه‌ریزان و پژوهشگران. جهت تصمیم‌گیری‌های راهبردی، تهیه گزارش‌های مختلف، توصیف وضعیت جاری می‌تواند مفید باشد. داده کاوی یا استخراج و کشف سریع و دقیق اطلاعات با ارزش و پنهان از این پایگاه داده‌ها از جمله اموری است که هر کشور، سازمان و شرکتی به منظور توسعه علمی، فنی و اقتصادی خود به آن نیاز دارد.

در کشور ما نیز سازمان‌ها، شرکت‌ها و مؤسسات دولتی و خصوصی از جمله سازمان تامین اجتماعی به طور فزاینده اقدام به ایجاد یا خرید نرم افزارهای پایگاه داده‌ها و مکانیزه کردن سیستم‌های اطلاعات خود هستند، همچنین با توجه به فصول دهم و یازدهم قانون برنامه سوم توسعه در خصوص داد و ستدهای الکترونیکی و همچنین تأکید بر برخورداری کشور از فن‌آوری‌های جدید اطلاعات برای دستیابی آسان به اطلاعات داخلی و خارجی، دولت مکلف شده است امکانات لازم برای دستیابی آسان به اطلاعات، زمینه سازی برای اتصال کشور به شبکه‌های جهانی و ایجاد زیرساخت‌های ارتباطی و شاهراه‌های اطلاعاتی فراهم کند. واضح است این امر باعث ایجاد پایگاه‌های عظیم داده‌ها شده و ضرورت استفاده از داده کاوی را که فرآیند خودکار کشف دانش و اطلاعات از پایگاه‌های داده می‌باشد، بیش از پیش نمایان می‌سازد.

هدف از داده کاوی ایجاد مدل‌هایی برای تصمیم‌گیری است. این مدل‌ها رفتارهای آینده را براساس تحلیل‌های گذشته پیش‌بینی می‌کنند. به کاربردن داده کاوی به عنوان اهرمی برای آماده سازی داده‌ها و تکمیل قابلیت‌های انبار داده، بهترین موقعیت را برای به دست آوردن برتری‌های رقابتی یا خدمت‌رسانی بهتر به مشتری ایجاد می‌کند.

با توجه به اینکه سازمان تامین اجتماعی و دیگر سازمان‌های بیمه دیگر مانند سازمان بیمه خدمات درمانی، سازمان بازمشستگی کشوری و ... همواره با داده‌ها و اطلاعات بسیار زیادی در مورد سوابق بیمه شده‌ها، کارفرمایان، اطلاعات درمانی بیمه شده‌ها، اطلاعات پرسنل، منابع مادی و ... روبرو هستند و در اکثر مواقع این داده‌ها می‌تواند حامل اطلاعات و الگوهای باارزشی باشند، لذا یکی از مهمترین کاربردهای داده‌کاوی در اینگونه سازمانها است.

امروز بانک‌های اطلاعاتی وسیعی از ویژگی‌های بیمه شده‌ها در سازمان تامین اجتماعی موجود است که اطلاعات مربوط به ویژگی‌های خانوادگی، تحصیلی و ... را شامل می‌شود. یافتن الگوها و دانش نهفته در این اطلاعات به تصمیم‌گیرندگان در این زمینه کمک شایانی خواهد کرد. استفاده از تکنیک‌های

پیشرفته داده‌کاوی مانند خوشه‌بندی، طبقه‌بندی، و ... می‌تواند در طبقه‌بندی کارفرمایان و کارگاه‌ها، مراکز درمانی، یافتن الگوهای خاص و با ارزش در مورد مراکز درمانی موفق، یافتن استراتژی توزیع تجهیزات درمانی، منابع مالی یافتن نقاط بحرانی در مدیریت مالی و موارد دیگر کاربرد داشته باشد. داده‌کاوی می‌تواند برای پاسخ دادن به یک سوال خاص مربوط به بیمه شده و نیز برای کشف روندهای عمومی که به تصمیم‌گیری کمک می‌کنند، استفاده شود. برای مثال سوال می‌تواند چنین باشد: امکان اینکه بیمه شده بعد از ده سال از شروع بیمه از کار افتاده شود چقدر است؟ یا میزان اعتبار مورد انتظار برای مستمری‌بگیران سازمان یا اعتبار لازم برای درمان در سال آینده چقدر است؟ درک الگوی استفاده کلی از خدمات درمانی یا تحلیل درخواستهای برای طی ۵ سال گذشته نیز همگی مثالهایی از کشف روندهای عمومی اند.

هم‌اینک در شرکت مشاور مدیریت و خدمات ماشینی تامین دانش و شناخت کافی بر روی شیوه‌ها و ابزارهای داده‌کاوی فراهم شده است و در این زمینه و یافتن الگوهای متفاوت بر مبنای اطلاعات جمع‌آوری شده فعالیت می‌شود.

نوشته ویوک.ان. پاتکار (V. N. Patkar)

ترجمه مریم صراف زاده و افسانه حاضری

Email: mmsarraaf@yahoo.com

دانشجویان دکتری سیستمهای اطلاعاتی- ملیورن استرالیا

چکیده:

کتابخانه ها و موسسات آموزشی با مشکل مدیریت کارآمد بار سنگین داده ها که دائما نیز در حال افزایش است روبرو می باشند. نرم افزارهای کامپیوتری بکار گرفته شده برای این منظور، غالبا فقط برای پرس و جوهای معمولی و پشتیبانی از مسائل مدیریتی و برنامه ریزی کوتاه مدت اداری جوابگو هستند. در حالیکه در عمق درون این حجم داده ها، الگوها و روابط بسیار جالبی میان پارامترهای مختلف بصورت پنهان باقی میماند. داده کاوی یکی از پیشرفتهای اخیر در حوزه کامپیوتر برای اکتشاف عمیق داده هاست. داده کاوی از اطلاعات پنهانی که برای برنامه ریزیهای استراتژیک و طولانی مدت میتواند حیاتی باشد پرده برداری میکند. تبیین مشخصه های اساسی فراینده داده کاوی و کشف کاربردهای ممکن آن در کتابداری و موسسات دانشگاهی اهداف اصلی این مقاله را شکل میدهند.

مقدمه

در دنیای بشدت رقابتی امروز، اطلاعات بعنوان یکی از فاکتورهای تولیدی مهم پدیدار شده است. در نتیجه تلاش برای استخراج اطلاعات از داده ها توجه بسیاری از افراد دخیل در صنعت اطلاعات و حوزه های وابسته را به خود جلب نموده است. حجم بالای داده های دائما در حال رشد در همه حوزه ها و نیز تنوع آنها به شکل داده متنی، اعداد، گرافیکها، نقشه ها، عکسها، تصاویر ماهواره ای و عکسهای گرفته شده با اشعه ایکس نمایانگر پیچیدگی کار تبدیل داده ها به اطلاعات است. علاوه بر این، تفاوت وسیع در فرآیندهای تولید داده مثل روش آنالوگ مبتنی بر کاغذ و روش دیجیتال مبتنی بر کامپیوتر، مزید بر علت شده است. استراتژیها و فنون متعددی برای گردآوری، ذخیره، سازماندهی و مدیریت کارآمد داده های موجود و رسیدن به نتایج معنی دار بکار گرفته شده اند. بعلاوه، عملکرد مناسب ابرداده [۱] که داده ای درباره داده است در عمل عالی بنظر میرسد.

پیشرفتهای حاصله در علم اطلاع رسانی و تکنولوژی اطلاعات، فنون و ابزارهای جدیدی برای غلبه بر رشد مستمر و تنوع بانکهای اطلاعاتی تامین می کنند. این پیشرفتهای هم در بعد سخت افزاری و هم نرم افزاری حاصل شده اند. ریزپردازنده های سریع، ابزارهای ذخیره داده های انبوه پیوسته و غیر پیوسته، اسکنرها، چاپگرها و دیگر ابزارهای جانبی نمایانگر پیشرفتهای حوزه سخت افزار هستند. پیشرفتهای حاصل در نظامهای مدیریت بانک اطلاعات در طی چهار دهه گذشته نمایانگر تلاشهای بخش نرم افزاری است. این تلاشها در بخش نرم افزار را میتوان بعنوان یک حرکت

پیشرونده از ایجاد یک بانک اطلاعات ساده تا شبکه ها و بانکهای اطلاعاتی رابطه ای و سلسله مراتبی برای پاسخگویی به نیاز روزافزون سازماندهی و بازیابی اطلاعات ملاحظه نمود. بدین منظور در هر دوره، نظامهای مدیریت بانک اطلاعاتی [۲] مناسب سازگار با نرم افزار سیستم عامل و سخت افزار رایج گسترش یافته اند. در این رابطه میتوان از محصولاتمانند، Unify, Dbase-IV, Sybase, Oracle و غیره نام برد.

داده کاوی یکی از پیشرفتهای اخیر در راستای فن آوریهای مدیریت داده هاست. داده کاوی مجموعه ای از فنون است که به شخص امکان میدهد تا ورای داده پردازی معمولی حرکت کند و به استخراج اطلاعاتی که در انبوه داده ها مخفی و یا پنهان است کمک می کند. انگیزه برای گسترش داده کاوی بطور عمده از دنیای تجارت در دهه ۱۹۹۰ پدید آمد. مثلا داده کاوی در حوزه بازاریابی، بدلیل پیوستگی غیرقابل انتظاری که بین پروفایل یک مشتری و الگوی خرید او ایجاد میکند اهمیتی خاص دارد. (Barry and Linoff, 1997)

تحلیل رکوردهای حجیم نگهداری سخت افزارهای صنعتی، داده های هواشناسی و دیدن کانالهای تلویزیونی از دیگر کاربردهای آن است. در حوزه مدیریت کتابخانه کاربرد داده کاوی بعنوان فرایند مآخذ کاوی [۳] نامگذاری شده است. این مقاله به کاربردهای داده کاوی در مدیریت کتابخانه ها و موسسات آموزشی می پردازد. در ابتدا به چند سیستم سازماندهی داده ها که ارتباط نزدیکی به داده کاوی دارند می پردازد؛ سپس عناصر داده ای توصیف میشوند و در پایان چگونگی بکارگیری داده کاوی در کتابخانه ها و موسسات آموزشی مورد بحث قرار گرفته و مسائل عملی مرتبط در نظر گرفته می شوند.

پیشرفت در تکنولوژیهای داده پردازی

سازمانهای بزرگ و چند-مکانه مثل بانکها، دفاتر هواپیمایی و فروشگاههای زنجیره ای با حجم زیادی از داده ها که ناشی از عملکرد روزانه آنهاست روبرو هستند. بطور سنتی چنین داده هایی به دو دسته تقسیم شده اند:

۱. رکوردهای اصلی [۴]

۲. رکوردهای عملیاتی [۵]

فرض بر این است که رکوردهای اصلی حاوی اطلاعات پایه هستند که معمولا چندان تغییر نمی کنند در حالیکه رکوردهای عملیاتی با توجه به طبیعت عملیات تجاری حتی بطور ساعتی تغییر خواهند کرد.

سیستمهای مدیریت پایگاه داده [۶] مناسب برای پیوند دادن این دو مجموعه اطلاعاتی و تهیه گزارشهای استاندارد جهت کنترل فعالیتها گسترش یافتند. سیستم اطلاعات مدیریت رایج برای پشتیبانی عملیات و سرویس دهی به چند کاربر در سطوح مختلف سازمان مبتنی بر این نظریه است.

بمنظور کمک به تصمیم گیری راهبردی، نظریه تاسیس بانک اطلاعات رکوردهای اصلی به نظریه سازماندهی دیتا مارت [۷] و انبار داده ها [۸] تغییر یافت. استخراج اطلاعات از رکوردهای عملیاتی یا پایگاههای اطلاعات عملیاتی و سازماندهی آن برای تحلیل استاندارد یا زمانی فلسفه اولیه و اصولی چنین پیشرفتهایی است. گرچه، دیتا مارت و انبار داده ها از نظر هدف و ساختار با هم

دیتامارت

دیتا مارت اغلب کوچک است و بر یک موضوع یا دپارتمان خاص متمرکز است. بنابراین پاسخگوی یک نیاز داخلی است. طرح بانک اطلاعات برای یک دیتامارت حول ساختار اتصال ستاره ای ساخته شده است که بهینه برای نیازهای کاربران دپارتمان است. دیتامارت معمولا با ابزارهای کامپیوتری که انعطاف پذیری تحلیل را تامین میکنند اما ممکن است برای سازماندهی حجم بالای داده ها مناسب نباشند؛ نیرومند میشود. رکوردهای ذخیره شده در دیتامارتهای بخوبی نمایه شده اند. یک دیتامارت در صورتیکه داده ها را از منابع داده ای بسیار سازماندهی شده مثل انبار داده ها بگیرد؛ دیتامارت وابسته نامیده میشود. مسلما دیتامارتهای وابسته از لحاظ ساختاری و معماری منطقی هستند. منبع دیتامارتهای وابسته تکنولوژی بانک اطلاعات دپارتمانی است. دیتامارتهای مستقل ثابت نیستند و از لحاظ معماری بسیار با هم متفاوتند. این مساله هنگام یکپارچه سازی دیتامارتهای مستقل، مشکل ایجاد میکند. بنابراین با یکپارچه سازی ساده دیتامارتهای یک انبار داده ایجاد نخواهد شد. دیتامارت اساسا برای اهداف تاکتیکی طراحی شده است و هدفش تامین یک نیاز تجاری فوری است.

انبار داده ها

یک انبار داده کاملا " متفاوت از دیتامارت است. سازماندهی انبارهای داده بگونه ایست که کلیه موضوعات حول فعالیتهای کاری سازمان را می پوشاند. انبار داده نمایانگر یک تسهیلات مرکزی است. برخلاف دیتامارت که در آن داده ها به شکل خلاصه تر و متراکم تر وجود دارند، یک انبار داده ، داده ها را در یک سطح نامتراکم ذخیره می کند. ساختار داده ها در یک انبار داده یک ساختار لزوما" هنجار شده است. بدین معنی که ساختار و محتوای داده ها در انبار داده منعکس کننده ویژگیهای دپارتمانهای عضو نیست. داده ها در انبار داده از نظر حجم و شکل کاملا" متفاوت از داده ها در دیتامارت هستند. دیتامارت ممکن است شامل حجم زیادی از داده های قدیمی و گذشته نگر باشد. داده ها در انبار داده اغلب بصورت نسبتا" سبک نمایه میشوند. (به بیان دیگر در عمق کمتر). انبار داده برای اهداف برنامه ریزی بلندمدت و راهبردی طراحی میشوند. در نتیجه انبار داده برخلاف سیستم عملیات که کاربرمدار است متمرکز بر اقلام است. ساختار یک انبار داده مشخصات زیر را نشان میدهد:

وابستگی به زمان:

رکوردها بر اساس یک برچسب زمانی نگهداری میشوند. وابستگی زمانی حاصل در ایجاد صفحات زمانی مفید است که درک ترتیب زمانی وقایع را تسهیل میکند.

غیر فرار بودن[۹]:

رکوردهای داده در انبار داده ها هرگز بطور مستقیم روزآمد نمیشوند. برای هر تغییری در ابتدا داده های عملیاتی روزآمد میشوند و سپس بگونه ای مقتضی به انبار داده منتقل میشوند. این مساله

ثبات داده ها را برای استفاده های وسیعتر تضمین میکند.

تمرکز موضوعی:

داده ها از بانکهای اطلاعاتی عملیاتی بصورت گزینشی به انبار داده منتقل میشوند. این استراتژی به ایجاد یک انبار داده بر اساس یک مطلب یا موضوع خاص کمک میکند و بنابراین کاوش انبار داده ها برای پرس و جوهای موضوعی با سرعت بیشتری انجام میشود.

یکپارچگی:

داده ها بگونه ای کامل سازماندهی شده اند تا با حذف موارد تکراری و چند عنوانه یکپارچگی رکوردها حفظ شود؛ به ایجاد ارجاع های متقابل کارآمد بین رکوردها کمک نموده و ارجاع دهی را تسهیل نماید.

واضح است که انبار داده اساساً برای پرس و جوهای پشتیبان تصمیم گیری ساخته شده است. بر این اساس سازماندهی و عملیات انبار داده چنان طراحی شده اند تا نیازهای اطلاعاتی روزمره یا معمولی را پاسخگو باشند. بدلیل حجم بسیار بالای چنین پایگاه اطلاعاتی یک سیستم کامپیوتری پیشرفته برای عملیات انبارسازی داده ها لازم است. همچنین یک بانک اطلاعات مجزا شامل ابرداده که مشخصه هایی نظیر نوع، فرمت، مکان و پدیدآورندگان داده های ذخیره شده در یک انبار داده ها را توصیف میکند نیز برای کمک به کاربران و مدیران داده ها ساخته میشود. مشخص شد که انبار داده بدلیل اندازه و تنوعش، اگر مبتکرانه پردازش شود میتواند به تولید اطلاعاتی منجر شود که در وهله اول آشکار نیستند. با انتخاب متناسب داده ها، بکار گرفتن فنون مختلف غربال کردن و تفسیر زمینه ای [۱۰]، داده ذخیره شده میتواند منجر به کشف الگوها یا رابطه هایی شود که بینش نویی به تصمیم گیرنده دهد. این مساله نظریه توسعه عملیات داده کاوی را به موازات معدن کاوی بروز داد. ذکر این نکته لازم است که داده کاوی در اصل لزوماً نیاز به سازماندهی یک انبار داده ندارد. حال به داده کاوی می پردازیم.

عناصر داده کاوی

توصیف و کمک به پیش بینی دو کارکرد اصلی داده کاوی هستند. تحلیل داده مربوط به مشخصه های انتخابی متغیرها؛ از گذاشته و حال، و درک الگو مثالی از تحلیل توصیفی است. برآورد ارزش آینده یک متغیر و طرح ریزی کردن روند مثالی از توانایی پیشگویانه داده کاوی است. برای عملی شدن هر یک از دو کارکرد فوق الذکر داده کاوی، چند گام ابتدایی اما مهم باید اجرا شوند که از این قرارند:

۱. انتخاب داده ها

۲. پاک سازی داد ها

۳. غنی سازی داده ها

۴. کد گذاری داده ها

با دارا بودن هدف کلی در مطالعه، انتخاب مجموعه داده های اصلی برای تحلیل، اولین ضرورت است. رکوردهای لازم میتواند از انبار داده ها و یا بانک اطلاعاتی استخراج شود. این رکوردهای داده جمع آوری شده؛ اغلب از آنچه آلودگی داده ها نامگذاری شده است رنج می برند و

بنابراین لازم است پاکسازی شوند تا از یکدستی فرمت (شکلی) آنها اطمینان حاصل شود، موارد تکراری حذف شده و کنترل سازگاری دامنه بعمل آید. ممکن است داده های گردآوری شده از جنبه های خاصی ناقص یا ناکافی باشند. در این صورت داده های مشخصی باید گردآوری شوند تا بانک اطلاعات اصلی را تکمیل کنند. منابع مناسب برای این منظور باید شناسایی شوند. این فرایند مرحله غنی سازی داده ها را تکمیل میکند. یک سیستم کدگذاری مناسب معمولاً "جهت انتقال داده ها به فرم ساختار-بندی شده جدید؛ متناسب برای عملیات داده کاوی تعبیه میشود .

فنون داده کاوی

ممکن است متوجه شده باشید که فنون داده کاوی یک گروه نامتجانس را شکل میدهند چرا که هر تکنیکی که بتواند بینش جدیدی از داده ها را استخراج کند میتواند داده کاوی به حساب آید. برخی از ابزارهای رایج بکار گرفته شده تحت عنوان داده کاوی عبارتند از: (Adriaans and Zantinge, 2003)

ابزارهای پرس و جو [۱۱]: ابزارهای متداول زبان پرس و جوی ساختار-بندی شده [۱۲] در ابتدا برای انجام تحلیلهای اولیه بکار گرفته شدند که می تواند مسیرهایی برای تفحص بیشتر نشان دهد.

فنون آماری: مشخصات اصلی داده ها لازمست با کاربرد انواع مختلفی از تحلیلهای آماری شامل جدول بندی ساده [۱۳] و متقاطع [۱۴] داده ها و محاسبه پارامترهای آماری مهم بدست آید.

مصور سازی: با نمایش داده ها در قالب نمودارها و عکسها مانند نمودار پراکنندگی؛ گروه بندی داده ها در خوشه های متناسب تسهیل میشود. استنباط عمیق تر ممکن است با بکارگیری تکنیکهای گرافیکی پیشرفته حاصل شود.

پردازش تحلیلی پیوسته [۱۵]: از آنجا که مجموعه داده ها ممکن است روابط چندین بعدی داشته باشند، روشهای متعددی برای ترکیب کردن آنها وجود دارد. ابزارهای پردازش تحلیلی پیوسته به ذخیره چنین ترکیباتی کمک میکند و ابزارهای ابتدا-انتها [۱۶] پیوسته برای انجام پرس و جو ایجاد میکند. اما این ابزارها هیچ دانش جدیدی ایجاد نمی کنند.

یادگیری مبتنی بر مورد: این تکنیک مشخصات گروههای داده ها را تحلیل میکند و به پیش بینی هر نهاد واقع شده در همسایگی شان کمک میکند. الگوریتمهایی که استراتژی یادگیری تعاملی را برای کاوش در یک فضای چندین بعدی بکار میگیرند برای این منظور مفیدند.

درختان تصمیم گیری: این تکنیک بخشهای مختلف فهرست پاسخهای موفق داده شده مربوط به یک پرس و جو را بازبایی می کند و به این ترتیب به ارزیابی صحیح گزینه های مختلف کمک میکند.

قوانین وابستگی: اغلب مشاهده میشود که یک وابستگی نزدیک (مثبت یا منفی) بین مجموعه ای از داده های معین وجود دارد. بنابراین قوانین رسمی وابستگی برای تولید الگوهای جدید ساخته و بکار گرفته میشوند.

شبکه های عصبی : این یک الگوریتم یادگیری ماشینی است که عملکرد خودش را بر اساس کاربرد و ارزیابی نتایج بهبود می بخشد.

الگوریتم ژنتیکی: این هم تکنیک مفید دیگری برای پیش بینی هدف است. به این ترتیب که با یک گروه یا خوشه شروع میشود و رشدش در آینده را با حضور در برخی مراحل فرایند محاسبه احتمال جهش تصادفی؛ همانطور که در تکامل طبیعی فرض میشود طرح ریزی می نماید. این تکنیک به چند روش میتواند عملی شود. و ترکیب غیرقابل انتظار یا نادری را از عواملی که در حال وقوع بوده و مسیر منحنی طراحی داده ها را تغییر میدهند؛ منعکس میکند.

گام نهایی فرایند داده کاوی، گزارش دادن است. گزارش شامل تحلیل نتایج و کاربردهای پروژه، در صورت بکارگیری آنها، است. و متن مناسب، جداول و گرافیکها را در خود جای می دهد. بیشتر اوقات گزارش دهی یک فرایند تعاملی است که تصمیم گیرنده با داده ها در پایانه کامپیوتری بازی میکند و فرم چاپی برخی نتایج واسطه محتمل را برای عملیات فوری بدست می آورد.

داده کاوی در تولید چهار نوع دانش ذیل مفید است: (Fayyad et al., 1996)

- دانش سطحی (کاربردهای SQL)
- دانش چند وجهی (کاربردهای OLAP)
- دانش نهان (تشخیص الگو و کاربردهای الگوریتم یادگیری ماشینی)
- دانش عمیق (کاربردهای الگوریتم بهینه سازی داخلی)

نرم افزار:

از آنجا که داده کاوی با بانکهای اطلاعاتی بزرگ سروکار دارد، به گونه ای ایده ال با تکنولوژی خدمت گیر-خدمت گر [۱۷] بکار میرود. کاربردهای عمومی داده کاوی بیشتر شامل تقسیم کردن داده ها در خوشه های مقتضی، کدگذاریهای مناسب، کاوش برای الگوها و طراحی کردن با استفاده از فنون آماری و الگوریتمهای ژنتیکی است. تعداد زیادی از بسته های نرم افزاری واجد این جنبه های ابزارهای داده کاوی با درجات متفاوتی از جامعیت در دسترس هستند. برای مثال بسته های نرم افزاری که منحصرًا برای کاربردهای OLAP در دسترس هستند عبارتند از: Oracle OLAP, DB2 OLAP Server, CleverPath OLAP. نرم افزارهای آماری عمومی مثل SPSS, SAS, STATISTICA با امکاناتی برای داده کاوی و بسته های نرم افزاری اختصاصی داده کاوی مثل Weka, Insightful Miner3, Text Mining Software, Enterprise Data Mining software, PolyAnalyst 4.6 مفید هستند.

کاربردهای داده کاوی در کتابخانه ها و محیط های دانشگاهی

داده کاوی در ابتدا از حوزه تجارت برخاست اما کاربردهای آن در سایر حوزه هائی که به گردآوری حجم وسیعی از داده هائی می پردازند که دستخوش تغییرات پویا نیز می گردند؛ مفید شناخته شد. بخشهایی مثل بانکداری، تجارت الکترونیک، تجارت سهام، بیمارستان و هتل از این نمونه اند. انتظار میرود که استفاده از داده کاوی در بخش آموزش بطور عام امکانات جدید بسیاری ارائه دهد. برخی کاربردهای داده کاوی در کتابخانه ها و قسمت اداری آموزش در ذیل مورد بحث قرار گرفته اند.

مدیریت و خدمات کتابخانه

عملیات کتابداری بطور کلی شامل مدیریت مدارک، ارائه خدمات و امور اداره و نگهداری است. هر کدام از این کارکردها با انواع مختلفی از داده ها سروکار دارد و بطور جداگانه پردازش میشود. اگرچه، انجام تحلیل ترکیبی براین مجموعه های داده نیز میتواند افق تازه ای را بگشاید که به طرح خدمات جدید و تحول رویه ها و عملیات جاری کمک نماید. جدول یک برخی از کاربردهای ممکن داده کاوی را که میتواند در کتابداری مفید باشد ارائه میکند.

جدول یک- کاربردهای داده کاوی در کتابخانه ها

کاربرد متصور	بانک اطلاعاتی
برای تعیین نقاط قوت و ضعف مجموعه	گردآوری منابع
برای ایجاد رابطه بین خواننده، منابع کتابخانه و زمان مشخصی از سال	استفاده از مجموعه
برای تحلیل سفارشهای پاسخ داده شده و سفارشهای دریافت شده	امانت بین کتابخانه ای
برای پیش بینی روند بازگشت منابع	داده های بخش امانت
برای نشان دادن منابع مالی بکار گرفته شده	داده های هزینه

داده کاوی میتواند برای پاسخ دادن به یک سوال خاص مربوط به کتابخانه و نیز برای کشف روندهای عمومی که به تصمیم گیری کمک میکنند، استفاده شود. برای مثال سوال میتواند چنین باشد: امکان اینکه امانت گیرندگان منابع را یک هفته بعد از تاریخ عودت برگردانند تا نامه های یادآوری کمتری فرستاده شود چقدر است؟ یا میزان اشتراک مورد انتظار برای نشریات بین المللی انتخاب شده برای سال آینده چقدر است؟ درک الگوی استفاده کلی مجلات الکترونیکی یا تحلیل درخواستهای اعضا برای میکروفیلرها طی ۵ سال گذشته نیز همگی مثالهایی از کشف روندهای عمومی اند. دامنه تحلیل استنادی هم میتواند با استفاده از داده کاوی گسترش داده شود. در ارتباط با کتابخانه ها، وب کاوی حوزه دیگری از علاقمندی است. وب کاوی شامل محتوا کاوی وب، ساختار کاوی وب و استفاده کاوی وب با توجه به یک موضوع خاص است که در طراحی خدمات جدید مبتنی بر وب کمک خواهد کرد.

مدیریت موسسات دانشگاهی

اداره موسسات دانشگاهی کار پیچیده ای است. در این موسسات داتما" نیاز به درآمدزایی و خود-کارآمدی و کاهش وابستگی به بودجه دولتی احساس میشود. این مساله کنترل دائمی جنبه های مختلف هر فعالیت و پروژه را می طلبد. بانکهای اطلاعاتی برای چنین موسساتی مربوط به دانشجویان، دانشکده، اساتید و کارمندان، تعداد رشته ها و چند مورد دیگر است. ارزیابی تقاضا و وضعیت عرضه نقش مهمی بازی میکند. مرور بانکهای اطلاعاتی نمونه در جدول ۲ نمایانگر کاربردهای بالقوه داده کاویست.

جدول ۲- کاربردهای داده کاوی در موسسات دانشگاهی

کاربرد متصور	بانک اطلاعاتی
--------------	---------------

ثبت نام دانشگاهی	برای درک رابطه های جمعیت شناختی، اقتصادی و اجتماعی
کارایی دانشگاهی	برای ایجاد رابطه بین عوامل اقتصادی-اجتماعی و نمرات اخذ شده
بانک سوالات	برای تعیین میزان مفید بودن سیستم با استناد به نمرات امتحان
همکاری فکری	برای ارزیابی همکاری دانشکده با توجه به میزان استفاده از کتابخانه
انتشارات	برای پیدا کردن تأثیر انتشارات در تقاضا برای رشته ها
بازدید از وب سایت	برای تحلیل سوالات دریافت شده در وب سایت دانشگاه و کمک به ایجاد رشته های جدید دانشگاهی

کاربرد داده کاوی در دانشگاه ملی سنگاپور قابل ملاحظه است. در این دانشگاه از ابزارهای داده کاوی برای شناسایی و دسته بندی دانشجویانی که به کلاسهای پیش نیاز برای واحد درسی ارائه شده نیاز داشتند استفاده شد. (Kurian and John, 2005) علاوه بر آن، مسائلی مانند اختصاص بهتر منابع و نیروی انسانی، مدیریت روابط دانشجو و به تصویر کشیدن رفتار گروههای مختلف میتواند بوسیله ابزارهای داده کاوی انجام شود.

محدودیت ها

کاربرد داده کاوی با چند عامل محدود شده است. اولین مورد به سخت افزار و نرم افزار لازم و موقعیت بانک اطلاعاتی مربوط میشود. برای مثال در هند، داده های غیر مجتمع که برای کاربردهای داده کاوی لازم است ممکن است به فرم دیجیتالی در دسترس نباشد. در دسترس بودن نیروی انسانی ماهر در داده کاوی نیز مسأله مهم دیگری است. محرمانه بودن رکوردهای مراجعان ممکن است در نتیجه پردازش داده های مبتنی بر داده کاوی آسیب پذیر شود. کتابداران و مؤسسات آموزشی باید این مسأله را در نظر داشته باشند؛ چرا که در غیر اینصورت ممکن است گرفتار شکایات قانونی گردند.

محدودیت دیگر از ضعف ذاتی نهفته در ابزارهای نظری ناشی میگردد. ابزارهایی مانند یادگیری ماشینی و الگوریتمهای ژنتیکی بکار گرفته شده در فعالیتهای داده کاوی به مفاهیم و فنون منطق و آمار بستگی دارد. در این حد نتایج به روش مکانیکی تولید شده و بنابراین به یک بررسی دقیق نیاز دارند. اعتبار الگوهای بدست آمده به این طریق؛ باید آزمایش شود. چرا که در بسیاری موارد روابط علل و معلول مشتق شده؛ از برخی استدلالات غلط ذیل رنج میبرند. (Cannavo, 2003)

• علت دور

مثلاً" امکانات ضعیف خوابگاه باعث می شود دانشجویان نمرات پایینی کسب نمایند

• علت مجرد

مثلاً" بودجه محدود بر بازدهی پژوهشی دانشکده تأثیر می گذارد

• علائم در نظر گرفته شده برای این عوامل

مثلاً" مجموعه کتابخانه ممکن است افزایش نیابد چون تعداد خوانندگان مرتباً کاهش می یابد.

• سفسطه دسته بندی

مثلاً" مدرسان حقوق بسیار بالا دریافت می کنند و کل حقوقشان بالغ بر میلیونها میشود.

- سفسطه ترکیب

مثلاً اگر هر مدرس در دانشکده شایسته و واجد صلاحیت باشد کل دانشکده عملکرد بهتری خواهد داشت.

- سوگیری در انتخاب نمونه:

مثلاً استناد به یافته های یک پیمایش نمونه گیری شده از دانشجویان یک دانشکده که از خانواده های ثروتمند هستند و مخارج روزانه در خوابگاه برای هر دانشجو ۱۰۰ دلار است. از آنجایی که مطالعه الگوها و استخراج روابط میان رکوردها مستلزم کاربرد منطق قیاسی و استقرایی است فرد باید مراقب اشتباهاتی که عموماً رخ میدهد باشد. برای مثال بحثهای قیاسی یا استقرایی، تا زمانیکه وضعیت درست بودن فرضیه آزمایش نشود چیزی درباره درست یا غلط بودن نتایجشان نمی گویند. طبیعتاً، نتایج تولید شده ماشینی ممکن است از چنین نقایصی رنج ببرند.

تذکرات نهایی

بکارگیری تکنولوژی اطلاعات توسط هر سازمان در عمل یک فرایند هموار نیست. کتابخانه یا مؤسسه دانشگاهی از این قضیه استثناء نیست. اما، تجربه نشان میدهد که یک برنامه نظام مند میتواند ظهور و نگهداری تکنولوژی اطلاعات در محیط کتابخانه را تسهیل کند. (Patkar and Iyer, 1990; Patkar, 2000, 2004) حتی کاربرد تکنولوژی های پیشرفته پردازش اطلاعات مثل سیستمهای خبره و سیستم اطلاعات جغرافیایی (جی.آی.اس) در کتابخانه گزارش شده است. (Patkar, 1999; Myers, 1992) با این پیش زمینه، کاربرد داده کاوی بوسیله کتابخانه ها و موسسات دانشگاهی، به شرط آماده سازی مناسب، بطور قابل توجهی عملی است. برای دانشگاهها، کالجها، مدارس و موسسات آموزش از راه دور که بانکهای اطلاعاتی عظیمی دارند، ابزارهای داده کاوی میتواند الگوها و روابطی را که خیلی عیان نیستند آشکار کند. این نتایج ممکن است به طراحی دوباره فرایندها و رویه های مرتبط منجر شود. تحلیلهای پشتیبانی شده توسط داده کاوی در کل موسسات و محیط ها میتواند مسائل متنوع مدیریت آموزشی؛ از جمله درک بهتر مشخصه های اقتصادی اجتماعی دانشجویان، مندرجات رشته ها و آموزش و پرورش و ساختار هزینه را مخاطب قرار دهد.

آنچه لازم است اینست که فراتر از عملکرد داده پردازش استاندارد قدم برداریم مخصوصاً کتابخانه ها و موسسات دانشگاهی که با انواع مختلفی از بانکهای اطلاعاتی سروکار دارند و به سطوح معقولی از کامپیوتری کردن و دیجیتالی کردن داده ها دست یافته اند. در یک نظر، ابزارهای داده کاوی نمایانگر پیشرفت در زنجیره تکنولوژی اطلاعات هستند. داده کاوی همچنین میتواند بعنوان بخشی از فرایند بزرگتر کشف دانش در بانکهای اطلاعاتی در محیط های مختلف در نظر گرفته شود. البته نباید چنین پنداشت که ابزارهایی مثل داده کاوی نیاز به مداخله انسانی را کاهش خواهد داد. همچنانکه در بالا نشان داده شد، ارزیابی و تعدیل نتایج بدست آمده بوسیله چنین ابزارهای خودکاری؛ به آزمایش نیاز دارد تا در برابر کاربردهای غلط محافظت شود. انتظار می رود داده کاوی در گسترش سازمان خودیادگیرنده مشارکت کند. کشف انتخابهای نوین با بهره گیری از داده کاوی اطمینان بخش بهترین کاربرد ممکن منابع موجود است. داده کاوی ماهیت چرخه مانند دارد. برای اینکه در پی کشف الگوها، سوالات بیشتری پدید خواهند آمد که دور بعدی فرایند را شکل میدهند. بهره برداری از تکنولوژیهای پیشرفته مثل داده کاوی مطمئناً برای

متخصصان کتابداری و مدیران موسسات آموزشی یک چالش دائمی خواهد بود ؛ چرا که آنها خلاقیت طلبند و برای نوآوری تلاش می کنند.

!Error

REFERENCES

1. Adriaans, P; Zantinge, D. *Data Mining. 9th Indian reprint*. Delhi: Pearson Education, 2003.
2. Barry, M; Linoff, G. *Data Mining Techniques for Marketing, Sales and Customer Support*. New York: John Wiley & Sons, 1997.
3. Cannavo, S. *Think to Win: The Power of Logic in Everyday Life*. Mumbai: Magna Publishing Co. Ltd. 2003.
4. Fayyad, U.M; Piatetsky-Shapiro G; Smyth, P; Uthurusamy, R. *Advances in Knowledge Discovery and Data Mining*. Cambridge MA: AAAI Press/MIT Press, 1996
5. Inmon, B. Data Mart does not Equal Data Warehouse. *DM Review*. May 1998. (http://www.dmreview.com/article_sub.cfm?articleId=608).
6. Kurian, J.C; John, B.M. *Mining the Education Domain*. http://www.getforme.com/previous2004/151004_miningtheeducationdomain.htm, (accessed on March 31, 2005).
7. Myers, J.E. Reference Expert: Developing a Computer Expert System for Library Reference Service. Paper presented at the *Conference on Advances in IT Applied to Libraries*, University of Puerto Rico, Rio Piedras Campus, 8 October 1992.
8. Nicholson, S. The Bibliomining Process: Data Warehousing and Data Mining for Library Decision-Making. *Information Technology and Libraries*. Vol.22(4); 2003; p4-9.
9. Patkar, V.N; Iyer, P.V. Determining Priorities for Computerisation in Libraries. *Annals of Library Science and Documentation*. Vol. 37(3); 1990; p110-116.
10. —. Application of GIS in Library Management. *Information Studies*. Vol. 5(2); 1999; p73-82.
11. —. Information Technology Adoption by Libraries: Opportunities and Barriers. *Information Studies*. Vol. 6(2); 2000; p125-140.
12. —. Adaptation to Information Technology Development by Libraries. *Information Studies*. Vol. 10(4); 2004; p637-648.

یادداشتها:

- [1] Metadata
- [2] DBMS
- [3] bibliomining
- [4] Master records
- [5] Transaction records

- [6] DBMS=Database Management System
- [7] Data Mart
- [8] Data warehouse
- [9] Non-volatility
- [10] Contextual interpretation
- [11] Query tools
- [12] Structured Query Language(SQL)
- [13] Simple tabulation
- [14] Cross tabulation
- [15] Online analytical processing (OALP)
- [16] Front-end tools
- [17] Client-server

* این مقاله ترجمه ای است از:

Patkar, V. N. (2005). "Data mining applications in library and academic institutions." *Information Studies* 11(3): 145-156.

داده کاوی در مدیریت ارتباط با مشتری
پرستو شاه سمندی
Parastoushahsamandi@yahoo.com

چکیده

شرکتهای امروزی از طریق تجزیه و تحلیل چرخه زندگی مشتری به افزایش ارزش مشتری دست یافته اند. ابزارها و فناوریهای انبار داده، داده کاوی و دیگر تکنیک های مدیریت ارتباط با مشتری، روشهایی هستند که فرصتهای جدیدی را برای تجارت فراهم کرده اند.

در واقع دیدگاه محصول محوری جای خود را به مشتری محوری داده است. بنابراین، با جمع آوری داده های مربوط به مشتری و تصمیم گیری براساس الگوهای استخراج شده از روابط پنهان میان داده ها به وسیله ابزار داده کاوی، می توان به خواسته مشتری محوری خود جامه عمل پوشاند. این مقاله مفاهیمی از مدیریت ارتباط با مشتری و یکی از عناصر آن - داده کاوی- را مورد بررسی قرار می دهد.

مقدمه

در سالهای اخیر فرهنگ تجارت به پیشرفتهایی نایل گشته است. مطابق با آن روابط اقتصادی مشتریان به شیوه های بنیادی و اساسی در حال تغییر است. شرکتها به منظور نظارت بر اینگونه تغییرات نیازمند آرایه راه حلها هستند. ظهور و پیدایش اینترنت در تغییر جهت مرکز توجه بازاریابی نقش بسزایی داشته است. چنانچه اطلاعات بر خط (ON LINE) بیشتر در دسترس قرار گیرد موجب آگاهی و هوشیاری بیشتر مشتریان می گردد. آنها در جریان تمام آنچه آرایه و پیشنهاد می شود قرار می گیرند و تقاضای بهترینها را دارند. برای از عهده برآمدن در چنین شرایطی باید سیستم هایی که بتواند به طور دقیق نسبت به مشتریان واکنش نشان دهد به کار رود. جمع آوری آمار مشتریان و داده های رفتاری آنها این هدف اصلی و دقیق را ممکن می سازد. این نوع هدفگیری به يك برنامه ریزی عالی هنگام ایجاد يك رقابت سخت و به مشخص کردن مشتریان بالقوه هنگام عرضه محصولات جدید کمک می کند.

داده کاوی

امروزه با حجم عظیمی از داده ها روبرو هستیم. برای استفاده از آنها به ابزارهای کشف دانش نیاز داریم. داده کاوی به عنوان يك توانایی پیشرفته در تحلیل داده و کشف دانش مورد استفاده قرار می گیرد. داده کاوی در علوم (ستاره شناسی،...) در تجارت (تبلیغات، مدیریت ارتباط با مشتری،...) در وب (موتورهای جستجو،...) در مسایل دولتی (فعالیتها ضد تروریستی،...) کاربرد دارد. (۱) عبارت داده کاوی شباهت به استخراج زغال سنگ و طلا دارد. داده کاوی نیز اطلاعات را که در انبارهای داده مدفون شده است، استخراج می کند. (۲)

در واقع هدف از داده کاوی ایجاد مدل هایی برای تصمیم گیری است. این مدلها رفتارهای آینده را براساس تحلیلهای گذشته پیش بینی می کنند. به کاربردن داده کاوی به عنوان اهمی برای آماده سازی داده ها و تکمیل قابلیتهای انبار داده (DATA WAREHOUSE)، بهترین موقعیت را برای به دست آوردن برتریهای رقابتی ایجاد می کند.

سیستم های بانک داده (BASE DATA) ، نقشی کلیدی در سیستم های مدیریت و انبار داده، بازی می کنند. یک سیستم بانک داده، شامل فایل های بانک داده و سیستم های مدیریت بانک داده است. (۱)

اغلب تجارت ها به تصمیم گیریهای استراتژیک و یا اتخاذ خط مشی های جدید برای خدمت رسانی بهتر به مشتریان نیاز دارند. به عنوان مثال فروشگاهها آرایش مغازه خود را برای ایجاد میل بیشتر به خرید مجدداً طراحی می کنند و یا خطوط هواپیمایی تسهیلات خاصی را برای مشتریان جهت پروازهای مکرر آنها در نظر می گیرند. این دو مثال به داده هایی در مورد رفتار مصرفی گذشته مشتریان برای تعیین الگوهای به وسیله داده کاوی، نیاز دارد. براساس این الگوها تصمیمات لازم اتخاذ می شود. در واقع ابزار داده کاوی، داده را می گیرد و یک تصویر از واقعیت به شکل مدل می سازد، این مدل روابط موجود در داده ها را شرح می دهد. (۲)

از نظر فرایندی فعالیتهای داده کاوی به سه طبقه بندی عمومی تقسیم می شوند: (۶)

اکتشاف : فرایند جستجو در یک بانک داده برای یافتن الگوهای پنهان، بدون داشتن یک فرضیه از پیش تعیین شده درباره اینکه این الگو ممکن است چه باشد.

مانند تحلیلهایی که برحسب کالاهای خریداری شده صورت می گیرد، اینگونه تحلیلهای سبدي نشانگر مواردیست که مشتری تمایل به خرید آنها دارند. این اطلاعات می تواند به بهبود موجودی، استراتژی طراحی، آرایش فروشگاه و تبلیغات منجر گردد.

مدل پیش بینی : فرایندی که الگوهای کشف شده از بانک داده را می گیرد و آنها را برای پیش بینی آینده به کار می برد.

مانند پیش بینی فروش در خرده فروشی، الگوهای کشف شده برای فروش به آنها کمک می کند تا تصمیماتی را در رابطه با موجودی اتخاذ کنند.

تحلیلهای دادگاهی : به فرایند به کارگیری الگوهای استخراج شده برای یافتن عوامل داده ای نامعقول و متناقض مربوط می شود.

مانند شناسایی و تشخیص کلاهبرداری در موسسات مالی. کلاهبرداری به میزان زیادی پرهزینه و زیان آور است، بانکها می توانند با تحلیل دادوستدهای جعلی گذشته الگوهای را برای تشخیص و کشف کلاهبرداری به دست آورند.

مدیریت ارتباط با مشتری

مدیریت ارتباط با مشتری یک فرایند تجاری است که تمام جوانب مشخصه های مشتری را آدرس دهی می کند، دانش مشتری را به وجود می آورد، روابط را با مشتری شکل می دهد و برداشت آنها را از محصولات یا خدمات سازمان ایجاد می کند. مدیریت ارتباط با مشتری توسط چهار عنصر از یک چارچوب ساده تعریف شده است: دانش، هدف، فروش و خدمت. (۳)

مدیریت ارتباط با مشتری با در نظر گرفتن اینکه چه محصولات یا خدماتی، به چه مشتریانی، در چه زمانی و از طریق چه کانالی عرضه شود، بهبود را در پی خواهد داشت. این مدیریت از اجزای مختلفی تشکیل شده است.

پیش از اینکه فرایند آن آغاز شود، شرکت باید اطلاعات مشتری را در اختیار داشته باشد. این اطلاعات می تواند از داده های داخلی مشتریان و یا از داده های منابع خارجی خریداری شده، به دست آید. برای داده های داخلی منابع مختلفی وجود دارد مانند پرسشنامه ها و بلاگ ها، سوابق کارت اعتباری و...

منابع داده خارجی یا بانکهای داده خریداری شده مانند آدرسها، شماره تلفن ها، پروفایل های بازدید از وب سایتها کلیدی برای به دست آوردن دانش بیشتری از مشتری است. (۳)

بیشتر شرکتها، بانکهای داده ای عظیمی شامل داده های بازاریابی، منابع انسانی و مالی را دارا هستند. بنابراین، سرمایه گذاری در زمینه انبار داده، یکی از اجزای حیاتی در استراتژی مدیریت ارتباط با مشتری است. (۴)

پس از تهیه و تخصیص منابع داده، سیستم مدیریت ارتباط با مشتری باید با به کارگیری ابزارهایی مانند داده کاوی، داده ها را تجزیه و تحلیل کند. اعم از اینکه شرکت تکنیک های آماری سنتی را به کار می برد یا یکی از ابزارهای نرم افزاری مانند داده کاوی را، کارشناسان نیاز به فهم داده های مشتری و روابط تجاری دارند. بنابراین، داشتن افرادی متخصص که این داده ها را با ابزارهای مربوطه استخراج و به صورت اطلاعات درآورند، مهم است.

چرخه زندگی مشتری

واژه چرخه زندگی مشتری به مراحل در ارتباط بین مشتری و تجارت بر می گردد و آگاهی نسبت به آن موجب سودآوری مشتری می شود. عموماً چهار مرحله در چرخه زندگی مشتری وجود دارد:

۱ - **مشتریهای بالقوه**: افرادی که هنوز مشتری نیستند ولی در هدف بازار قرار دارند؛

۲ - **مشتریهایی که عکس العمل نشان می دهند**: مشتریان بالقوه یا احتمالی که به یک محصول یا خدمت علاقه و واکنش نشان می دهند.

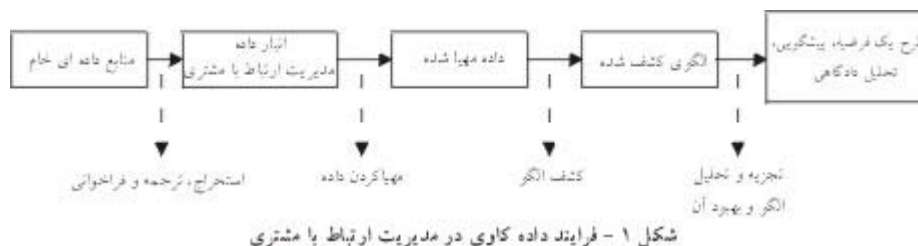
۳ - **مشتریهای بالفعل**: افرادی که در حال حاضر محصول یا خدمتی را به کار می برند.

۴ - **مشتریهای سابق**: اینگونه افراد مشتریان مناسبی نیستند چرا که مدت زیادی در هدف فروش قرار ندارند و یا خریدشان را به سمت محصولات رقیب برده اند. (۲)

فرایند داده کاوی در مدیریت ارتباط با مشتری

داده کاوی یکی از عناصر مدیریت ارتباط با مشتری است و می تواند به حرکت شرکتها به سمت مشتری محوری کمک کند.

فرایند داده کاوی در مدیریت ارتباط با مشتری به صورت زیر است. (شکل ۱)



داده های خام از منابع مختلفی جمع آوری می شوند و از طریق استخراج، ترجمه و فراوندهای فراخوانی به انبار داده این مدیریت وارد می شوند. در بخش مهیاسازی داده، داده ها از انبار خارج شده و به صورت يك فرمت مناسب برای داده کاوی در می آیند.

بخش کشف الگو شامل چهار لایه است:

- ۱ - سوالهای تجاری مانند توصیف مشتری، ۲ - کاربردها مانند امتیازدهی، پیش گوئی، ۳ - روشها مانند سری های زمانی، طبقه بندی، ۴ - الگوریتم ها.
- در این بخش روشهای داده کاوی با کاربرد مخصوص خود برای پاسخ به سوالهای تجاری که به ذهن می رسند، الگوریتم هایی را استخراج می کنند و از این الگوریتم ها برای ساخت الگو استفاده می شود.
- در بخش تجزیه و تحلیل الگو، الگوها به يك دانش مفید و قابل استفاده تبدیل می شوند و پس از بهبود آنها، الگوهایی که کارا محسوب می شوند در يك سیستم اجرایی به کار گرفته خواهند شد. (۱)

نتیجه گیری

رابطه مشتری با زمان تغییر می کند و چنانچه تجارت و مشتری درباره یکدیگر بیشتر بدانند این رابطه تکامل و رشد می یابد. چرخه زندگی مشتری چارچوب خوبی برای به کارگیری داده کاوی در مدیریت ارتباط با مشتری فراهم می کند. در بخش ورودی داده کاوی، چرخه زندگی مشتری می گوید چه اطلاعاتی در دسترس است و در بخش خروجی آن، چرخه زندگی می گوید چه چیزی احتمالاً جالب توجه است و چه تصمیماتی باید گرفته شود. داده کاوی می تواند سودآوری مشتری های بالقوه را که می توانند به مشتریان بالفعل تبدیل شوند، پیش بینی کند و اینکه تا چه مدت به صورت مشتریان وفادار خواهند ماند و چگونه احتمالاً ما را ترک خواهند کرد.

بعضی از مشتریان مرتباً مراجعاتشان را به شرکتها برای کسب مزیتهایی که طی رقابت میان آنها به وجود می آید، تغییر می دهند. در این صورت شرکتها می توانند هدفشان را روی مشتریانی متمرکز کنند که سودآوری بیشتری دارند.

بنابراین می توان از طریق داده کاوی ارزش مشتریان را تعیین، رفتار آینده آنها را پیش بینی و تصمیمات آگاهانه ای را در این رابطه اتخاذ کرد.

منابع:

۱ - STATE UNIVERSITY, 2003 NONE YE, "THE HAND BOOK OF DATA MINING", ARIZONA

CUSTOMER RELATIONSHIP CHRIS RYGIELSKI, "DATA MINING TECHNIQUES FOR - 2
.MANAGEMENT", TECHNOLOGY IN SOCIETY, 2002
.M., "THE CUSTOMER LIFECYCLES", INTELLIGENT ENTERPRISE, 1999 FREEMAN - 2
.CRM: EASIER SAID THAN DONE", INTELLIGENT ENTERPRISE, 1999" ,.HILL L - 2
GEMINI, "FOUR ELEMENTS OF CUSTOMER RELATIONSHIP CAP & IDC - 0
.PAPER MANAGEMENT", CAP GEMINI WHITE
INFORMATION DISCOVERY INC., "A CHARACTERIZATION OF DATA MINING - 1
[HTTP://WWW.DMREVIEW.COM/PORTAL](http://www.dmreview.com/portal) ,"TECHNOLOGIES AND PROCESSES

- Why Data Mining?
- What Is Data Mining?
- Potential Applications
- Data Mining Process
- Data Mining and Business Intelligence
- Interdisciplinary aspects of Data Mining
- Why Not Traditional Data Analysis?
- Multi- Dimensional view of Data Mining
- What Kinds of Data?
 - Data Mining tasks
- Architecture: Typical Data Mining System
- History of Data Mining
- Conferences and Journals on Data Mining
- Where to Find References
- Data Preprocessing
- Data Warehousing and OLAP Technology
- Mining Frequent Patterns, Association and Correlations
- Classification and Predication (LDA, QDA, Decision trees, Neural Networks, Regression, Clustering, etc)
- Numerical Examples for Classification and Prediction using SPSS Clementine ●Applications and Trends in Data Mining

مرجع: کتاب داده کاوی - دکتر جمال شهرابی

داده کاوی

مقاله زیر را من از اینترنت گرفتم، متأسفانه منبع دقیق آن را یادم نیست. در هر صورت امیدوارم مفید باشد.

داده کاوی

داده کاوی فرایندی تحلیلی است که برای کاوش داده ها (معمولاً حجم عظیمی از داده ها - در زمینه های کسب و کار و بازار) صورت می گیرد و یافته ها با به کارگیری الگوهای، احراز اعتبار می شوند. هدف اصلی داده کاوی پیش بینی است. فرایند داده کاوی شامل سه مرحله می باشد: ۱. کاوش اولیه ۲. ساخت مدل یا شناسایی الگو با کمک احراز اعتبار/ تایید و ۳. بهره برداری.

مرحله ۱: کاوش. معمولاً این مرحله با آماده سازی داده ها صورت می گیرد که ممکن است شامل پاک سازی داده ها، تبدیل داده ها و انتخاب زیرمجموعه هایی از رکوردها با حجم عظیمی از متغیرها (فیلدها) باشد. سپس با توجه به ماهیت مساله تحلیلی، این مرحله به مدل های پیش بینی ساده یا مدل های آماری و گرافیکی برای شناسایی متغیرهای مورد نظر و تعیین پیچیدگی مدل ها برای استفاده در مرحله بعدی نیاز دارد.

مرحله ۲: ساخت و احراز اعتبار مدل. این مرحله به بررسی مدل های مختلف و گزینش بهترین مدل با توجه به کارایی پیش بینی آن می پردازد. شاید این مرحله ساده به نظر برسد، اما اینطور نیست. تکنیک های متعددی برای رسیدن به این هدف توسعه یافتند. و "ارزیابی رقابتی مدل ها" نام گرفتند. بدین منظور مدل های مختلف برای مجموعه داده های یکسان به کار می روند تا کار □ پی شان با هم مقایسه شود، سپس مدلی که بهترین کارایی را داشته باشد، انتخاب می شود. این تکنیک ها عبارتند از:

Meta-learning, Boosting, Stacking

مرحله ۳: بهره برداری. آخرین مرحله مدلی را که در مرحله قبل انتخاب شده است، در داده های جدید به کار می گیرند تا پیش بینی های خروجی های مورد انتظار را تولید نماید. داده کاوی به عنوان ابزار مدیریت □ اطلاعات برای تصمیم گیری، عمومیت یافته است. اخیراً، توسعه تکنیک های تحلیلی جدید در این زمینه مورد توجه قرار گرفته است (مثلاً Classification Trees)، اما هنوز داده کاوی مبتنی بر اصول آماری نظیر (Exploratory Data Analysis (EDA) می باشد.

با این وجود تفاوت عمده ای بین داده کاوی و EDA وجود دارد. داده کاوی بیشتر به برنامه های کاربردی گرایش دارد تا ماهیت اصلی پدیده. به عبارتی داده کاوی کمتر با شناسایی روابط بین متغیرها سروکار دارد.

مفاهیم اساسی در داده کاوی

Bagging: این مفهوم برای ترکیب رده بندی های پیش بینی شده از چند مدل به کار می رود. فرض کنید که قصد داریم مدلی برای رده بندی پیش بینی بسازیم و مجموعه داده های مورد نظرمان کوچک است. شمامی توانید نمونه هایی (با جایگزینی) را از مجموعه داده ها انتخاب و برای نمونه های حاصل از درخت رده بندی (مثلاً C&RT و CHAID) استفاده نمایید. به طور کلی برای نمونه های مختلف به درخت های متفاوتی خواهید رسید. سپس برای پیش بینی با کمک درخت های متفاوت به دست آمده از نمونه ها، یک رای گیری ساده انجام دهید. رده بندی نهایی، رده بندی ای خواهد بود که درخت های مختلف آنرا پیش بینی کرده اند.

Boosting: این مفهوم برای تولید مدل های چندگانه (برای پیش بینی یا رده بندی) به کار می رود.

Boosting نیز از روش C&RT یا CHAID استفاده و ترتیبی از classifier ها را تولید خواهد کرد.

Meta-Learning: این مفهوم برای ترکیب پیش بینی های حاصل از چند مدل به کار می رود. و هنگامی که انواع مدل های موجود در پروژه خیلی متفاوت هستند، کاربرد دارد. فرض کنید که پروژه داده کاوی شما شامل Tree classifier ها نظیر C&RT و CHAID، تحلیل خطی و شبکه های عصبی است. هر یک از کامپیوترها، رده بندی هایی را برای نمونه های پیش بینی کرده اند. تجربه نشان می دهد که ترکیب پیش بینی های چند روش دقیق تر از پیش بینی های هر یک از روشهاست. پیش بینی های حاصل از چند classifier را می توان به عنوان ورودی meta-linear مورد استفاده قرار داد. meta-linear، پیش بینی ها را ترکیب می کند تا بهترین رده بندی پیش بینی شده حاصل شود.

امروزه **اطلاعات و فن آوری اطلاعات** آنچنان بر جنبه های مختلف زندگی سیاسی، اجتماعی، فرهنگی و اقتصادی افراد و اجتماعات تاثیر گذاشته است که نمی توان از آن غفلت نمود. تاثیر اطلاعات و فن آوری اطلاعات هم از جهت ایجاد فرصتهای طلایی جدید و هم از نظر چالشهای سازمانی قابل مطالعه است. بدون شك جهانی شدن اقتصاد، دسترسی به آن جوامع صنعتی و پیشرفته را به جامعه اطلاعاتی و اقتصادی آنها را از اقتصاد صنعتی به اقتصاد مبتنی بر خدمات اطلاعاتی و دانش تبدیل کرده است. فرآیند تولید و مدیریت شرکتهای و سازمانها به شدت متحول شده اند و سازمانها و شرکتهای مجازی به وجود آمده اند. فن آوری اطلاعات بسیاری از محدودیتهای زمانی و جغرافیایی را از بین بده است و افراد و سازمانها بدون توجه به مرزهای جغرافیایی می توانند با هم ارتباط برقرار کنند و به مبادله اطلاعات بپردازند. سیستمهای اطلاعاتی، ارتباط دوربرد و شبکه ها تحولات و دگرگونیهای اساسی در تولید کالا و خدمات، بازرگانی، بازاریابی، فروش، تبلیغ و مدیریت دسترسی به اطلاعات به وجود آورده اند. انقلاب اطلاعات که از قرن بیستم شروع شده است، تغییرات شگرفی را در نیروی کار به وجود آورده است و نیز حرفه ها و خدمات جدیدی ظهور کرده اند که به تولید، توزیع یا کار با دانش و اطلاعات وابسته اند، به طوری که بیش از ۶۰ درصد تولید ناخالص ملی و حدود ۵۵ درصد نیروی کارگری آمریکا به اطلاعات و دانش وابسته است (صفحه ۳ همین کتاب).

در تعریف اقتصاد به عنوان يك علم توافقی نظر کلی وجود ندارد. عده ای آن را مطالعه و بررسی تولید، توزیع، و مصرف ثروت در يك جامعه انسانی دانسته اند و گروهی دیگر آن را علم استفاده از منابع کمیاب به منظور نیل به هدف رفاه اجتماعی برای جامعه بشری. در مورد اینکه اقتصاد به عنوان يك علم مطرح است یا هنر نیز بین صاحب نظران اختلاف وجود دارد. عده ای معتقدند هنگامی که عملکردهای اقتصادی با استفاده از ریاضیات توجیه گردند. اقتصاد به عنوان يك علم مطرح است و بعضی دیگر معتقدند زمانی که از فلسفه و منطق برای توجیه مسایل اقتصادی استفاده کنند، اقتصاد به عنوان يك هنر مطرح خواهد بود. هر تعریفی را برای اقتصاد بپذیریم، مباحث اقتصادی مربوط به هر کالا یا خدمات به طور عمده شامل سرمایه گذاری، تولید، قیمت گذاری، بازاریابی، تبلیغ، توزیع و فروش مصرف درآمد، شرکتهای، اشتغال و سایر مواردی است که معمولاً در متون اقتصادی مطرح می گردند. بر این اساس اقتصاد اطلاعات مطالعه و بررسی تولید، توزیع، بازاریابی، قیمت گذاری، فروش، مصرف و کلیه درآمدهایی است که به طور مستقیم یا غیر مستقیم از طریق تولید، انتشار، فروش، ذخیره، پردازش، و دسترسی به اطلاعات حاصل می شود. اطلاعات و دسترسی به آن چنان ارزش و اهمیتی دارد که یکی از ارکان اساسی توسعه و محور تصمیم گیری در تمامی سطوح سیاسی، اقتصادی، فرهنگی و اجتماعی است و به همین دلیل هر سال سرمایه های هنگفتی برای تحقیق، تولید، انتشار، فراهم آوری، سازماندهی، ذخیره، بازاریابی، اشاعه و دسترسی به آن صرف یم گردد و به عنوان کالایی با ارزش خرید و فروش می گردد و بسیاری از خدمات و صنایع پیچیده برای ذخیره، پردازش، انتقال و دسترسی به آن به وجود آمده اند و بسیاری از کشورهای پیشرفته بریا جذب این بازار برنامه ریزی دقیق و سرمایه گذاری هنگفتی را انجام داده اند. اینترنت که پدیده ای جدید در فن آوری اطلاعات می باشد، بر اساس بررسی دانشگاه تگزاس، حدود ۳۰۰ میلیارد دلار درآمد در سال ۱۹۹۸ برای آمریکائیان داشته است. در حالی که اقتصاد جهانی سالانه به طور متوسط ۲/۸ درصد و در آمریکا ۲/۸ درصد رشد دارد، اقتصاد اینترنت به تنهایی در آمریکا به طور متوسط ۱۷۴/۵ درصد از سال ۱۹۹۵ تا ۱۹۹۸ رشد داشته است و در این

مدت کوتاه از صنایع قدیمی و قدرتمندی چون انرژی (۲۲۲ میلیارد دلار درآمد سالانه)، اتومبیل (۲۵۰ میلیارد دلار) و ارتباطات راه دور (۲۷۰ میلیارد دلار) در آمریکا پیشی گرفته است. در حالی که کشور آرژانتین از لحاظ درآمد ناخالص ملی در رتبه نوزدهم در جهان می باشد، درآمد حاصل از اینترنت در آمریکا به تنهایی بیشتر از درآمد ناخالص ملی کشور آرژانتین می باشد. اطلاعات، فناوریهای اطلاعات و ارتباطات به گونه ای توسعه و اهمیت یافتند که این عصر را عصر اطلاعات، ارتباطات و رایانه خوانده اند و خدمات و مشاغل زیادی را به خود جذب کرده است. برای مثال: بر اساس مطالعه و پیش بینی دانشگاه تگزاس تا سال ۱۹۹۸ تنها در بخش فن آوری اینترنت ۱/۲۰۰/۰۰۰ نفر شاغل بوده اند که عمدتاً شغل‌های جدیدی هستند. در انتشار و توزیع یک کتاب افرادی چون نویسنده، مترجم، ناشر، ویراستار، حروفچین، نمونه خوان، صفحه آرا، گرافیست، لیتوگراف، چاپچی، صحاف، مراکز پخش، کارگزاران و کتابفروشیها به طور مستقیم در ارتباطند. همچنین صنایع زیادی چون کاغذ، چاپ، سخت افزارها و نرم افزارهای رایانه ای نیز مستقیم و یا غیر مستقیم به آن مربوط می شوند. همچنین حرفه ها و مشاغل جدید دیگری چون فروشگاههای الکترونیکی کتاب (مانند آمازون)، تولید کنندگان و توزیع کنندگان کتابهای الکترونیکی نیز جزئی از اقتصاد کتاب محسوب می شوند که خود بخشی از اقتصاد اطلاعات می باشد. کشور آلمان در سال ۱۹۹۶، صرفنظر از سایر منابع اطلاعاتی، تعداد ۷۱۵۱۵ عنوان کتاب چاپ اول منتشر شده است. این بدان معنی است که نویسندگان، ناشران، و مشاغل و خدمات زیادی در این کشور به فعالیت مشغول بوده اند، به طوریکه تنها ارزش پشت جلد کتابها ۹ میلیارد دلار یا معادل ۹۰۰ میلیارد تومان میباشد. آمار منتشر شده از سوی انجمن ناشران آمریکا حاکی از فروش ۲۳/۰۲۳/۳۰۰/۰۰۰ دلار کتاب در سال ۱۹۹۸ می باشد. بر اساس آمار منتشر شده سال ۱۹۹۸ شرکت میکروسافت، که تنها به تولید نرم افزارهای رایانه ای می پردازد، ۲۵۱۸۲ نفر را به طور تمام در اختیار داشته است و در آمد خالص این شرکت در همان سال ۱۱/۳۶۰/۰۰۰/۰۰۰ دلار بوده است که نسبت به سال پیش ۵۴ درصد رشد داشته است. سفارش و تهیه، سازماندهی، حفاظت، نگهداری، ذخیره، پردازش و اشاعه اطلاعات، موجب تاسیس نهادها و شرکتهایی چون کتابخانه، مراکز اطلاعات و اطلاع رسانی، بانکهای اطلاعاتی، شرکتهای تولید انواع نرم افزار و سخت افزار و خدمات رایانه ای و بسیاری دیگر از این قبیل شده است. شبکه های اطلاعاتی عظیمی چون اینترنت به شدت در حلا دگرگون ساختن ارتباطات فرهنگی، اجتماعی، اقتصادی و ... است و ارزش و اهمیت اطلاعات و دسترسی به آن را دو چندان نموده است. مشاغل و خدمات جدیدی را به وجود آورده و مرزهای مکانی و زمانی را در دسترسی به اطلاعات از بین برده است. کشورهای پیشرفته با استفاده از فن آوری سنجش از راه دور و ماهواره ای به بسیاری از اطلاعات اقتصادی و نظامی سایر کشورها دسترسی دارند و با استفاده از این اطلاعات، که بسیار گران قیمت نیز هستند، روابط اقتصادی و نظامی خود را با دیگر کشورها تنظیم می کنند.

با استفاده وسیع از فن آوریهای جدید، اطلاعات به عنوان مواد خام به شکلهای مختلف مورد داد و ستد قرار می گیرد. می توانید دوره و شماره ای خاص از یک مجله را تصور کنید که گاهی به صورت چاپی و زمانی به شکلهای مختلف الکترونیکی (بر روی صفحه آرزان، صفحه فشرده یا سی.دی، دی.وی.دی. و اینترنت) مورد معامله قرار گیرد. مراکز چکیده نویسی و تولید پایگاههای اطلاعاتی مقالات همین مجله را به صورت دیگری بسته بندی و به مشتریان خود عرضه می نمایند. میزبانها اطلاعات آن را به صورت پیوسته و از طریق شبکه های رایانه ای در اختیار کاربران خود قرار می دهند. و به ازای هر بار بازیابی مقاله ای خاص هزینه ای دریافت می نمایند. هزینه دریافت این

خدمات ممکن است بر اساس میزان رکورد بازیابی شده، زمان اتصال و یا حتی سرعت انتقال اطلاعات باشد. متن کامل مقالات مجله مزبور زمانی دیگر توسط مراکز خدمات تحویل مدرک بارها فروخته می شود. با توجه به پیشرفت‌های اخیر در تولید، ذخیره، بازیابی و انتقال اطلاعات، بسیاری از ناشران و مراکز خدمات تحویل مدرک نوعی دیگر از تجزیه اطلاعات مربوط به همان مجله را رونق داده اند. تمامی خدمات مربوط به اشتراک، سفارش، دریافت، مسایل مالی و .. یک مجله از طریق اینترنت نیز صورت می پذیرد و کاربر می توان به صورت‌های مختلف اطلاعات یک مجله را از طریق اینترنت و از کانالها و شرکت‌های متفاوت دریافت نماید: اطلاعاتی چون فهرست مندرجات، چکیده مقالات، خدمات آگاهی رسانی جاری، اشاعه اطلاعات گزیده، دریافت متن کامل مجله، دریافت متن کامل یک مقاله خاص و غیره. هر کدام از این خدمات دارای نرخها و سیستم‌های قیمت گذاری جداگانه ای هستند.

به این ترتیب نه تنها اطلاعات را به عنوان یکی از عوامل تولید در کنار عواملی چون کار، زمین و سرمایه قرار می دهند که در تمامی عرصه های فرهنگی، اجتماعی، اقتصادی و سیاسی به عنوان عنصری اساسی و تاثیر گذار رخ می نماید، بلکه خود موجب رونق صنایع اطلاعاتی، ارتباطی و سایر صنایع مرتبط با اطلاعات نیز شده است: به طوری که بدون حضور اطلاعات بسیاری از این صنایع ارزش وجود خود را از دست خواهند داد.

به طور خلاصه، برای تولید، تهیه، سازماندهی، ذخیره، پردازش و دسترسی به اطلاعات فعالیتها، خدمات و مشاغل زیادی به طور مستقیم یا غیر مستقیم در ارتباطند که بخشی از اقتصاد اطلاعات محسوب می شوند. بریا مثال آن دسته از رایانه، سخت افزارها، نرم افزارها و مشاغلی که بریا تولید و توزیع کتاب به کار رفته می شوند، جزئی از اقتصاد اطلاعات محسوب می شوند. علیرغم آنکه اقتصاد اطلاعات در کشورهای صنعتی و پیشرفته به خوبی از اهمیت برخوردار است و این کشورها برنامه ریزیهای وسیعی را برای کسب درآمد از این رهگذر شروع کرده اند، ولی در کشورهای جهان سوم و در حال توسعه از جمله ایران اهمیت آن به خوبی درک نشده است و یا هنوز زیر ساخت‌های لازم برای توسعه این بخش از اقتصاد وجود ندارد. در اینجا سعی می شود در بخش‌های مختلف به طور مختصر به چشم اندازهای اقتصادی اطلاعات اشاره گردد.

فن آوری اطلاعات

تکنولوژی اطلاع رسانی / اطلاعاتی مجموعه ابزارها، ماشینها، دانش فنی، روشها و مهارتهای استفاده از آنها در تولید، داد و گرفت، پردازش، انباشت، بازیافت، جابه جایی، انتقال و مصرف اطلاعات است، از ساده ترین تا پیچیده ترین، و از ولتیتترین تا پیشرفته ترین مراحل اطلاعاتی. با این تعریف تمامی سخت افزارها مانند رایانه و لوازم جانبی آن و همچنین ابزارهای ارتباط دوربرد، شبکه های اطلاع رسانی، اینترنت، و نیز بسیار از نرم افزارها دیگر که برای ذخیره، پردازش، آماده سازی، بازیابی و مصرف اطلاعات به کار می روند در این رده قرار می گیرند. همانگونه که قبلاً نیز اشاره شد این بخش از صنعت اطلاعات به سرعت در حال گسترش و دگرگونی است و میلیاردها دلار درآمد هر ساله عاید شرکت‌هایی می شود که در این بخش فعالیت دارند. هر چند تولید سخت افزارهای اطلاع رسانی در ایران به کندی در حال پیشرفت است و در مقایسه با کشورهای پیشرفته فاصله زیادی را شاهد هستیم، ولی مدیریت سخت افزار تولید نرم افزار بخشی از فن آوری اطلاعات است که با برنامه ریزی دقیق می توان از طریق آن درآمدهای زیادی را نصیب شرکتها و کشور کرد. تولید و فروش این نرم افزارها می توان در سطح ملی و بین المللی باشد. در سطح ملی نرم افزارهای

خوب موجب استفاده بهینه از سخت افزارهاي موجود براي دسترسي به اهداف سازماني و فردي
نيز مي گردد، كه برآ آن هزينه هاي زيادي صرف شده است . براي مثال : توليد يك نرم افزار ذخير و
بازيابي اطلاعات براي كتابخانه ها موجب صرفه جويي در هزينه، سرعت دسترسي به اطلاعات ،
افزايش كيفي دسترسي به اطلاعات ، صرفه جويي در نيروي انساني، مديريت مناسبتر بر مجموعه
و بسياري از مزايي ديگر مي شود . كافي است بار ديگر اشاره كنيم كه در آمد خالص شركت
مايكروسافت كه فقط به توليد نرم افزار مي پردازد در سال ۱۹۹۸ معادل ۱۱/۳۶۰/۰۰۰/۰۰۰ دلار بود
ه است و نسبت به سال قبل ۵۴ درصد رشد داشته است. و يا انترنت در سال ۱۹۹۸ ،
۲۰۰/۰۰۰/۰۰۰/۰۰۰ دلار درآمد براي آمريکاييان داشته است و از سال ۱۹۹۵ دسدي معادل ۱۷۴/۵
درصد داشته است.

اشتغال

هر چند رشد صنعت اطلاعات موجب کاهش نياز به نيروي انساني در بسياري از بخشهاي توليد و
خدماتي شده است ، ولي خود موجب ايجاد بسيار از شغلهاي جديد نيز شده است . براي مثال در
اواخر دهه ۱۹۸۰ در کشور آفريقاي جنوبي در بخش اطلاع رساني ۲/۰۰۰/۰۰۰ نفر شاغل بودند . و
يا در سال ۱۹۹۸ در ايالات متحده آمريکا ۱/۲۰۰/۰۰۰ نفر فقط در بخش اينترنت شاغل بودند. اگر در
نظر بگيريم به طور متوسط حداقل دو نفر (با توجه به تيراژ در ايران) به طور مستقيم و غير مستقيم
براي توليد ، توزيع و مصرف يك كتاب درگير باشند، افزايش عناوين و تيراژ كتاب و ساير منابع
اطلاعاتي مي تواند در ايجاد شغل جديد بسيار موثر بشاد . بنا بر اين رشد توليد كتاب و ساير منابع
اطلاعاتي نه تنها فعاليت فرهنگي و علمي است كه موجب ارتقاي سطح آگاهي و دانش مردم و
كاهش بزهكاريهاي اجتماعي مي شود ، بلكه فعاليتي اقتصادي نيز محسوب مي شود . صنايع
كاغذ سازي: فن آوري اطلاعات : تاليف و ترجمه ، حروفچيني ، نمونه خواني، ويراستاري ادبي و
علمي، صفحه آرايي، فيلم و زينك، صحافي، چاپ، ناشران، كتابفروشي، توزيع عمده كتاب و سيار
منابع اطلاعاتي چاپي و الكترونيكي : كتابخانه ها: توليد ، اطلاعات و اطلاع رساني است و جزئي از
اقتصاد اطلاعات مي توانند محسوب شوند.

توليد توزيع و فروش اطلاعات

اطلاعات منبع اصلي تصميم گيري است و آنچه ان اهميتي استراتژيك دارد كه آن را قدرت مي دانند.
تصميم گيري درست و به موقع و برنامه ريزي کوتاه مدت و بلند مدت در مسائل اقتصادي، سياسي،
فرهنگي، اجتماعي و ساير موارد مربوط به يك سازمان، شركت ، و يا در سطح ملي و بين المللي
به اطلاعات مناسب و درست بستگي دارد و عدم دسترسي به آن خسارتهاي اقتصادي و مالي
فروان و در بعضي مواقع جبران ناپذيري با به جاب خواهد گذاشت. توليد و انتشار اطلاعات علمي -
فني ، اقتصادي، آماري و يا ساير اطلاعات مورد نياز و دسترسي به آن جزئي جدائي ناپذير از حيات
اقتصادي يك کشور مي باشد. به طوري كه شرکتهای موفق و بین المللي برای رقابت در بازار و تولید
محصول جديد و يا بهينه سازي آن ميلياردها دلار صرف تحقيق و پژوهش مي نمايند. براي مثال
صنايع داروسازي در آلمان سالانه بيش از ۵۰/۰۰۰/۰۰۰/۰۰۰ دلار براي تحقيق و پژوهش هزينه مي
كنند . اطلاعات نوعي مواد خام است كه در كالاهاي ساخته شده موجود بوده و توليد كالا به وجود
آن بستگي دارد. كالاهاي ساخته شده از بسياري جهات اطلاعات منجمد محسوب مي شوند. در
جهان امروز با توجه به حجم اطلاعات منتشر شده ، بدون سازماندهي و پردازش اطلاعات

دسترسی به بسیاری از اطلاعات عملاً ناممکن است. به همین دلیل سازماندهی، ذخیره، پردازش و تولید پایگاههای اطلاعاتی و اشاعه آن بین متقاضیان اطلاعات، حرفه ها و مشاغل زیادی را به خود اختصاص داده است و از این راه بسیاری از سازمانها و شرکتهای انتفاعی و غیر انتفاعی درآمدهای هنگفتی را نصیب خود می سازند. برای مثال کتابخانه ملی بریتانیا در سال ۱۹۹۸ روزانه بیش از ۱۴۰۰۰ مقاله به افراد متقاضی در سراسر دنیا ارسال کرده است و از راه درآمد سرشاری را نصیب خود کرده است. این خدمات بخشی از فروش سالیانه کتابخانه ملی بریتانیا است و شامل فروش نرم افزار، پایگاههای اطلاعاتی، خدمات مرجع و سایر خدمات جنبی دیگر نمی شود. نکته بسیار قابل تامل در فروش اطلاعات، در مقایسه با سایر کالاها، این است که یک قلم از اطلاعات پس از یک بار فروش مجدداً نیز قابل فروش است. برای مثال یک مقاله ممکن است دهها بار فروخته شود چکیده آن به صورتی دیگر به فروش برسد در قالب فهرست مندرجات، خدمات آکادمی رسانی جاری و اشاعه گزینشی اطلاعات به فروش برسد؛ و ... به کارگیری رایانه و شبکه های اطلاع رسانی به خصوص اینترنت تولید، سازماندهی و دسترسی به اطلاعات را متحول نموده است و مشاغل و شرکتهای جدیدی را درگیر نموده است. ناشران، مراکز خدمات تحویل مدرک و میزبانان پایگاههای اطلاعاتی از این طریق به تولید، بازیابی، فرورش و انتقال پاره های مختلف اطلاعاتی به مشتریان خود می پردازند و افراد متقاضی اطلاعات حتی در منازل خود نیز می توانند به اطلاعات مورد نیاز دسترسی پیدا کنند. آنها برای هر جزء از اطلاعات و خدماتی که ارائه می نمایند هزینه دریافت می کنند و مشتریان نیز با توجه به کیفیت و سرعت دسترسی به اطلاعات موجود در این سیستمها به صورت روزانه و در بعضی مواقع حتی دقیقه ای روزآمد می شوند و بدین ترتیب جاذبه های فراوانی برای دسترسی به بسیاری از اطلاعات به وجود آورده اند. با توجه به حجم زیاد اطلاعات، واحد اطلاعات از یک کتاب، مجله و به طور کلی از یک واحد کلی تر به جزئی از یک کتاب، مجله و حتی قسمتی از یک مقاله و واحدی جزئی تر تبدیل شده است و پایگاههای اطلاعاتی زیادی بر این اساس تولید و به بازار عرضه شده است. بسیاری از محققان و جویندگان اطلاعات به جای مراجعه به قفسه های کتابخانه و یا سایر واحدهای نگهداری اطلاعات ابتدا به این پایگاهها مراجعه و پس از انتخاب اولیه در این پایگاهها بریا تهیه متن کامل مطلب مورد نیاز به کتابخانه یا سایر واحدهای نگهداری اطلاعات مراجعه می کنند. بعضی از این پایگاههای اطلاعاتی استفاده کننده را مستقیماً به سایتی که مقاله در آنجا موجود است متصل می کنند و وی می توان مطلب مورد نظر خود را مطالعه، نسخه برداری و یا سفارش نماید. بسته بندی مجدد و تولید اطلاعات سفارشی و ارائه سریع آن به متقاضیان، فعالیت دیگری است که در حال توسعه است. تولید و توزیع اطلاعات به صورت الکترونیکی به خصوص از طریق اینترنت و نیز ارتباط مستقیم محققان با یکدیگر و مبادله مستقیم و سریع اطلاعات و همچنین تشکیل گروههای مباحثه و گروههای خبری و سایر روشهای دسترسی به اطلاعات، رفتارهای اطلاع یابی محققان و دانشمندان را به شدت تحت تاثیر قرار داده است. و ناشران را به چاش فراخوانده است. بسیاری از شرکتها و سازمانها معاملات تجاری، خرید و فروش و حتی پرداختهای مالی مربوطه را از طریق این نظامها و شبکه های اطلاع رسانی انجام می دهند. تجارت الکترونیکی چشم انداز جدیدی است که بسیاری از شرکتهای معتبر را به سوی خود جلب نموده است.

به طور خلاصه در عصری که آن را با نامهای گوناگون عصر اطلاعات، عصر کامپیوتر، عصر ارتباطات و عصر ماهواره می نامند، یک یز به عنوان نقطه اشتراك تمامی این نامگذاریها وجود دارد و آن این است که قطعات اطلاعات در پیکره این ابزارها مبادله می گردد. بدون حضور اطلاعات این ابزارها

ارزش وجود خود را از دست خواهند داد . تولید و دسترسی به موقع به اطلاعات سیاسی، اقتصادی ، علمی - فنی و سایر اطلاعات مورد نیاز موجب تصمیم گیری خردمندانه است و قدرت محسوب می شود به همین دلیل سرمایه های هنگفتی برای تولید، ذخیره، پردازش و دسترسی به اطلاعات صرف می گردد و صنایع اطلاعاتی و اقتصاد اطلاعات با نرخ رشدی شگفت آور در حال دگرگون کردن موازنه های اقتصادی است . صنایع پیچیده و ظریف بیش از پیش به اطلاعات و فن آوری اطلاعات وابسته است. اطلاعات نوعی مواد خام است که در کالاهای ساخته شده موجود بوده و تولید کالا به وجود آن بستگی دارد . کالاهای ساخته شده از بسیاری جهات اطلاعات منجمد محسوب می شوند.

حمید محسنی

عضو هیات علمی مرکز اطلاعات و مدارك علمی وزارت كشاورى

موسسه پژوهشی داده پردازان گیتا - دبیرخانه دائمی کنفرانس داده کاوی

به شرکت کنندگان در هر دوره (مقدماتی ، متوسط ، پیشرفته) که در تمامی کارگاههای آموزشی آن دوره شرکت نموده اند گواهینامه معتبر اعطا می گردد .

به شرکت کنندگانی که دارای حداقل یکی از شرایط زیر باشند حداکثر ۱۵% تخفیف تعلق می گیرد .

- شرکت کنندگان در هر سه دوره مقدماتی ، متوسط و پیشرفته

- ارائه کارت دانشجویی

- شرکت گروهی بیش از ۵ نفر

دوره مقدماتی (E)

پیش نیاز	کد کارگاه	عنوان کارگاه
ندارد	E1	مروری بر دانش داده کاوی (1) Data Mining Concepts (1)
E1	E2	مروری بر دانش داده کاوی (2) Data Mining Concepts (2)
ندارد	E3	فرایند و معماری داده کاوی Architecture and Data Mining Process
ندارد	E4	زمینه کاری و فعالیت داده کاوی Data Mining Tasks
ندارد	E5	الگوریتم ها و تکنیک های داده کاوی Data Mining Techniques and Algorithms
ندارد	E6	پروژه ها و کاربرد های داده کاوی Data Mining Projects and Applications

مدرس : جناب آقای دکتر جمال شهرابی

هزینه ثبت نام کل دوره : ۵۰۰۰۰۰ ریال

هزینه به ازای هر کارگاه : ۱۰۰۰۰۰ ریال

تاریخ برگزاری : ۲۲ تیرماه ۱۳۸۷

مکان : دانشگاه صنعتی امیرکبیر

دوره متوسط (I)

پیش نیاز	کد کارگاه	عنوان کارگاه
E	I1	قوانین وابستگی Data Mining Task: Association Rules
E	I2	خوشه بندی Data Mining Task: Clustering
E	I3	کلاس بندی Data Mining Task: Classification
E	I4	مدل های پیش بینی پیشرفته Advanced Data Mining Task: Prediction Models
E	I5	داده کاوی زمان محور / مکان محور Data Mining Task: Spatial and Temporal Data Mining
E	I6	متن کاوی / وب کاوی Data Mining Task: Text Mining / Web Mining

مدرس : جناب آقای دکتر جمال شهرابی

هزینه ثبت نام : ۸۰۰۰۰۰ ریال

هزینه به ازای هر کارگاه : ۱۵۰۰۰۰ ریال

تاریخ برگزاری : ۲۳ تیرماه ۱۳۸۷

مکان : دانشگاه صنعتی امیرکبیر

دوره پیشرفته (A)

پیش نیاز	کد کارگاه	عنوان کارگاه
E	A۱	مروری بر داده کاوی در SQL Server ۲۰۰۵ Over View of SQL Server ۲۰۰۵ Data Mining
E,I۱,A۱	A۲	تحلیل سبد بازار در SQL Server ۲۰۰۵(DMX) با استفاده از قواعد وابستگی Market Basket DMX Tutorial with Association Rules in SQL Server ۲۰۰۵
E,I۲,A۱	A۳	بخش بندی بازار در SQL Server ۲۰۰۵(DMX) با استفاده از خوشه بندی Segmentation DMX with Clustering in SQL Server ۲۰۰۵
E,I۳,A۱	A۴	کلاس بندی و پیش بینی در SQL Server ۲۰۰۵(DMX) با استفاده از درخت تصمیم Classification and Prediction DMX with Decision Tree in SQL Server ۲۰۰۵
E,A۱	A۵	ارزیابی صحت مدل ها با استفاده از نمودار صحت Evaluation with View Mining Accuracy Charts

مدرس : سرکار خانم مهندس ونوس شکورنیا

هزینه ثبت نام : ۱۰۰۰۰۰۰ ریال

تاریخ برگزاری : ۲۴ تیرماه ۱۳۸۷

مکان : دانشگاه صنعتی امیرکبیر

Data Mining Tasks

- Classification/partitioning
- Clustering
- Association
- Segmentation
- Regression
- Advanced prediction modeling
- Temporal data mining
- Spatial data mining
- Time series forecasting
- Deviation and outlier detection
- Explorative and visual data mining
- Web mining
- Text Mining
- Mining semi-structured data
- Content mining and pattern mining
- Multimedia mining (audio/video)
- Explorative and visual data mining
- Others

Data Mining Algorithms

- Clustering algorithms
- Genetic algorithms and categorization techniques
- Fuzzy logic and rough sets

Conference Topics

Data Mining Process

- Data preparation techniques
- Data reduction methods
- Data cleaning and preparation
- Feature selection and transformation
- Sampling and rebalancing
- Missing value imputation
- Model selection/assessment and comparison
- Model comparison
- Model interpretation
- Others

Data Mining Applications

- Engineering
- System and Manufacturing
- Industry and government
- System planning and management
- Urban planning and management
- Logistics/Traffic management
- Science and technology
- Education
- Business/Industrial

- Artificial neural networks
- Decision trees/rule learners
- Statistical methods
- Case based reasoning
- Link and sequence analysis
- Others

Data Mining Integration

- Mining large scale data
- Multidimensional data
- Distributed and grid based data mining
- Data visualization
- Knowledge Discovery in Databases (KDD)
- Data and knowledge representation
- Data warehousing
- OLAP integration
- Others

- Marketing
- Finance and financial services
- Insurance
- Social science
- Military/Security
- Bioinformatics/Medicine
- Biological sciences
- Risk analysis
- Emergency planning services
- Health, safety and environment (HSE)
- Others

کشف پول شویی و فساد مالی با روش‌های داده کاوی Data Mining

حالت اول: فرض کنید یک بزه‌کار اقتصادی بخواهد یک میلیون دلار را از طریق واردات کالا پول‌شویی کند. ابتدا لازم است به عنوان یک وارد کننده محلی، با یک صادرکننده خارجی (شایدخودش) تباری کند. سپس کارهای زیر را انجام دهد:

- ۱- صادرکننده خارجی ده هزار تیغ را بازای هر عدد یک دهم سنت می‌خرد. (جمعا ۱۰۰۰ دلار)
- ۲- صادرکننده خارجی ده هزار تیغ را به یک صادرکننده محلی به قیمت هر عدد ۱۰۰ دلار می‌فروشد. (جمعا ۱ میلیون دلار).
- ۳- واردکننده محلی ده هزار تیغ را به قیمت واقعی هزار دلار دریافت می‌کند و یک میلیون دلار به صادرکننده خارجی می‌پردازد.
- ۴- **نتیجه:** واردکننده محلی یک میلیون دلار را به یک کشور خارجی با هزینه هزار دلار انتقال (شست و شو) داده است. عمل فوق با اضافه صورت حساب کردن در واردات کالا صورت می‌گیرد.

حالت دوم: فرض کنید یک بزه‌کار اقتصادی بخواهد یک میلیون دلار را از طریق صادرات کالا پول‌شویی کند. ابتدا لازم است به عنوان یک صادر کننده محلی، با یک وارد کننده خارجی تباری کند. سپس کارهای زیر را انجام دهد:

- ۱- بزه‌کار محلی، با یک میلیون دلار، ۲۰۰ ساعت طلای تجملی به قیمت هریک پنج هزار دلار با پول نقد می‌خرد. (جمعا یک میلیون دلار).
 - ۲- صادرکننده محلی، ۲۰۰ ساعت طلا را به یک وارد کننده خارجی به قیمت هر عدد ۵ دلار می‌فروشد. (جمعا یک هزار دلار).
 - ۳- واردکننده خارجی، ۲۰۰ ساعت طلا را دریافت کرده و یک هزار دلار برای صادر کننده محلی صورت حساب می‌کند.
 - ۴- صادر کننده خارجی، ساعت‌های طلا را در بازار به قیمت هریک پنج هزار دلار می‌فروشد. (جمعا یک میلیون دلار).
- نتیجه:** صادر کننده محلی یک میلیون دلار به یک کشور خارجی با هزینه ۱۰۰۰ دلار انتقال (شست و شو) داده است. عمل فوق با کم صورت حساب کردن در صادرات کالا صورت می‌گیرد.

استفاده از تجارت جهانی برای انتقال پول سیاه از یک کشور به کشور دیگر، یکی از روش‌های قدیمی برای فرار از حسابرسی دولتی است. این کار از طریق اضافه صورت حساب کردن واردات یا کم صورت حساب کردن صادرات انجام می‌شود. البته برعکس هر دو عمل فوق نیز امکان دارد. سازمان‌های اقتصادی و اطلاعاتی معمولا کارهایی برای کشف پول‌شویی از طریق درهای جلو (موسسات مالی و حسابرسی) انجام داده و از درهای پشت ساختمان (تجارت بین الملل) مقداری غافل می‌مانند. البته برعکس حالت فوق نیز در ایران رایج است.

زیاد صورت حساب شدن واردات سه خلاف را ممکن است در بر داشته باشد:

- ۱- خلاف گمرکی.
 - ۲- فرار از مالیات.
 - ۳- پول‌شویی.
- این پول‌ها ممکن است در اختیار سازمانی مثل القاعده قرار گیرد. **عملیات پول‌شویی از این راه‌ها و انتخاب بهترین راه با مطالعه شرایط محیطی هر کشور صورت می‌گیرد.** بعنوان مثال در آمریکا چون کنترل کمتری روی صادرات کالا می‌شود، بهترین راه، از طریق صادرات کالا به دیگر کشورها است.

روش داده‌کاوی (Mining Data):

تحلیل اطلاعات واردات و صادرات کالا در وزارت بازرگانی آمریکا، ممیزی مالیاتی و بانک اطلاعاتی کالاهای تجاری نشان داده است که با روش‌های داده‌کاوی می‌توان به سرنخ‌هایی دست پیدا کرد. با داشتن قیمت تقریبی یک واحد کالا در زمان مشخص می‌توان کالاهای صادر و وارد شده را بررسی و انحراف‌های غیرعادی و زیاد را نشان کرده و ممیزی نمود. این کار برای ۱۶۳۹۰ کالای وارداتی و ۸۵۶۸ کالای صادراتی در سال ۲۰۰۱ و برای ۲۳۰ کشور که با آمریکا مرز مرز بازرگانی داشتند انجام شد. تمام واردات و صادرات با حد بالا و پائین مقایسه و ثبت گردید. مقادیر دلاری و تعداد موارد مشکوک برای هر کشور تجمیع شد. کل پول منتقل شده به خارج از آمریکا در سال ۲۰۰۱ مبلغ ۱۵۶ میلیارد دلار بود. جدول ۱ نمونه‌ای از موارد فوق است.

کالا	کشور	قیمت واحد
دستمال کاغذی خشکبار چیچی تیغ	چین بلژیک ژاپن انگلیس	۴۱۲۱ دلار هر هزار گرم ۲۴۲۶ دلار هر هزار گرم ۴۸۹۶ دلار هر عدد ۱۱۳ دلار هر عدد
کشورهای مظنون به وجود نیروهای القاعده در آنها		
حوله پنبه ای آبینه تیغ تلمبه دستی میل لنگ	پاکستان اندونزی مصر مالزی عربستان	۱۵۴ دلار هر عدد ۱۶۵ دلار هر سانتی مترمربع ۲۳ دلار هر واحد ۵۰۰۰ دلار هر واحد ۱۵۲۰۰ دلار هر واحد

جدول ۱ قیمت‌های غیرعادی واردات آمریکا

مقدار محاسبه شده که احتمال دارد توسط القاعده جایجا شده باشد حدود ۲۷ / ۴ میلیارد دلار برای ۲۵ کشور در فهرست مظنونین می‌باشد. برای پنج کشور بالای فهرست حدود ۶۵ / ۳ میلیارد دلار است. با توجه به مقدار کل ۱۵۶ میلیارد دلار پولشویی که از این راه محاسبه شده و مقدار ۲۷ / ۴ میلیارد دلار که القاعده مظنون به پولشویی است، حدود ۳ درصد برای القاعده و **۹۷ درصد برای قاتل‌ها و بزه‌کاران اقتصادی آمریکا و جهان** می‌شود. جدول ۲ نمونه‌ایی از صادرات غیرعادی از آمریکا به دیگر کشورها می‌باشد.

کالا	کشور	قیمت واحد
الماس تزئینی بالابر بلدوز خودکششی پروژکتور ویدیو موشک و سکوی پرتاب	هند جامائیکا کلمبیا برزیل اسرائیل	۱۲ دلار هر قیراط ۲۸۴ دلار هر دستگاه ۱۷۴۲ دلار هر دستگاه ۲۴ دلار هر دستگاه ۵۲ دلار هر واحد
کشورهای مظنون به وجود نیروهای القاعده در آنها		
مانیتور رنگی ویدیو مانیتور رنگی ویدیو کفش ورزشی ایزوتوپ رادیواکتیو	اندونزی پاکستان اردن مصر	۲۳ دلار هر دستگاه ۲۲ دلار هر دستگاه ۴ دهم دلار هر جفت ۱ صدم دلار هر واحد

جدول ۲ قیمت‌های غیرعادی صادرات آمریکا

نتیجه‌گیری : با توجه به روش‌ها و امکان داده کاوی، می‌توان از اطلاعات بانک‌های عامل و بانک مرکزی ، وزارت بازرگانی ، گمرکات کشور ، پایانه‌های حمل و نقل کالا، وزارت صنایع و معادن و دیگر سازمان‌های ذیربط استفاده نموده و موارد مشکوک را ممیزی کرد.
پیش نیاز فوق وجود راه‌ها و اطلاعات زیاد و کافی در سازمان های مربوطه است. اگر سامانه‌هایی در سازمان‌های مربوطه وجود داشته باشند که بتواند اطلاعات را بصورت "آن لاین" و فوری در اختیار بگذارند، با توجه به اعلام بازرگانان در مبادی حمل و نقل، صدور و ورود کالا، می‌توان قبل از اقدام به صادر یا وارد کردن کالا، آن را ممیزی و کشف نموده و مانع از پولشویی و فساد اقتصادی گردید. کارهای آماری و اطلاعاتی از این نوع یک علم تقریبی است، اما با امکانات رایانه‌ایی و روش‌های داده‌کاوی (Data Mining)، می‌توان کیفیت و دقت آنرا بیشتر کرد.

داده کاوی (Data Mining)

داده کاوی یکی از پیشرفتهای اخیر در راستای فن آوریهای مدیریت داده هاست. داده کاوی مجموعه ای از فنون است که به شخص امکان میدهد تا ورای داده پردازشی معمولی حرکت کند و به استخراج اطلاعاتی که در انبوه داده ها مخفی و یا پنهان است کمک می کند. انگیزه برای گسترش داده کاوی بطور عمده از دنیای تجارت در دهه ۱۹۹۰ پدید آمد. مثلاً داده کاوی در حوزه بازاریابی، بدلیل پیوستگی غیرقابل انتظاری که بین پروفایل یک مشتری و الگوی خرید او ایجاد میکند اهمیتی خاص دارد.

تحلیل رکوردهای حجیم نگهداری سخت افزارهای صنعتی، داده های هواشناسی و دیدن کانالهای تلویزیونی از دیگر کاربردهای آن است. در حوزه مدیریت کتابخانه کاربرد داده کاوی بعنوان فرایند مآخذ کاوی نامگذاری شده است. این مقاله به کاربردهای داده کاوی در مدیریت کتابخانه ها و موسسات آموزشی می پردازد. در ابتدا به چند سیستم سازماندهی داده ها که ارتباط نزدیکی به داده کاوی دارند می پردازد؛ سپس عناصر داده ای توصیف میشوند و در پایان چگونگی بکارگیری داده کاوی در کتابخانه ها و موسسات آموزشی مورد بحث قرار گرفته و مسائل عملی مرتبط در نظر گرفته می شوند.

مدیریت ذخیره سازی و دستیابی اطلاعات

داده های اطلاعاتی (Data) به عنوان یکی از منابع حیاتی سازمان شناخته می شود و بسیاری از سازمان ها با اطلاعات و دانش سازمانی خود مانند سایر دارایی های ارزشمندشان برخورد می کنند. نکته: داده اطلاعاتی (Data) به اطلاعات خام سازمان اطلاق می شود و اطلاعات (Information) به داده های پردازش شده. همچنین داده های پردازش شده پس از طبقه بندی و آنالیز به دانش سازمان (Knowledge) تبدیل می گردند.

کاوش های ماشینی در داده ها یا داده کاوی (Data mining) را باید یکی از سامانه های هوشمند (Intelligent systems) دانست. سامانه های هوشمند زیر شاخه ایست بزرگ و پرکاربرد از یادگیری ماشینی که خود زمینه ایست در هوش مصنوعی. زمینه علمی جدید و پهناور یادگیری ماشینی (که "کاوش های ماشینی در داده ها" بخشی ست بزرگ از زیر شاخه سامانه های هوشمند آن ست)، به واقع

همان امتداد و استمرار دانش کهن و همه جا گیر آمار است در جهت ماشینی کردن یادگیری، تعلّم، و سرانجام، دانش.

داده کاوی به عنوان مهمترین کاربرد Data Warehouse یا انباره های داده شناخته می شود. به وسیله داده کاوی های موجود مورد تحلیل قرار می گیرند تا روندهای احتمالی، ارتباطهای غیر محسوس و الگوهای مخفی داده ها از بین انبوه داده ها، شناسایی شوند.

از نظر فرایندی فعالیتهای داده کاوی به سه طبقه بندی عمومی تقسیم می شوند:

اکتشاف: فرایند جستجو در یک بانک داده برای یافتن الگوهای پنهان، بدون داشتن یک فرضیه از پیش تعیین شده درباره اینکه این الگو ممکن است چه باشد.

مانند تحلیلهایی که برحسب کالاهای خریداری شده صورت می گیرد، اینگونه تحلیلهای سبدي نشانگر مواردیست که مشتری تمایل به خرید آنها دارند. این اطلاعات می تواند به بهبود موجودی، استراتژی طراحی، آرایش فروشگاه و تبلیغات منجر گردد.

مدل پیش بینی: فرایندی که الگوهای کشف شده از بانک داده را می گیرد و آنها را برای پیش بینی آینده به کار می برد.

مانند پیش بینی فروش در خرده فروشی، الگوهای کشف شده برای فروش به آنها کمک می کند تا تصمیماتی را در رابطه با موجودی اتخاذ کنند.

تحلیلهای دادگاهی: به فرایند به کارگیری الگوهای استخراج شده برای یافتن عوامل داده ای نامعقول و متناقض مربوط می شود.

مانند شناسایی و تشخیص کلاهبرداری در موسسات مالی. کلاهبرداری به میزان زیادی پرهزینه و زیان آور است، بانکها می توانند با تحلیل دادوستدهای جعلی گذشته الگوهایی را برای تشخیص و کشف کلاهبرداری به دست آورند.

از نمایی دیگر، داده کاوی، بعنوان روشی در استخراج دانش از متون، یکی از موضوعات مهم در گستره ای از اعمال مدیریت اطلاعات است. در این میان آنچه از اهمیت فوق العاده ای برخوردار است آرایه راه کارهایی برای مواجهه با این حجم عظیم اطلاعاتی و استفاده بهینه از اطلاعات در جهت خلق دانش، تولید سینرجی و در نهایت افزایش خرد جمعی است.

در سالهای اخیر اهمیت متون به عنوان منابع با پتانسیل اطلاعاتی بسیار بالا به نحو گسترده ای مورد توجه قرار گرفته به طوری که کشف دانش از متون به عنوان یکی از مهمترین فعالیتهای محققین حوزه هوش مصنوعی و فناوری اطلاعات قرار گرفته است. تحقیقات بسیاری صورت گرفته اما محدوده فعالیت بقدری گسترده است که نیازمند توجه بیشتری می باشد.

ان نویسنده/نویسندگان: ترجمه و تالیف: مهندس عدرا قبادی
چکیده: جامعه مبتنی بر اطلاعات را می توان به عنوان جامعه ای تعریف نمود که بخش غالب اجتماع به جای کارهای فیزیکی در گیرکارهای فکری هستند. در چنین جامعه ای بیشترین توجه به فعالیتهای اطلاعاتی از قبیل: فراهم آوری، پردازش، تولید، ثبت، انتقال، اشاعه و مدیریت اطلاعات مبذول می گردد و بیشترین هزینه ها صرف فرایندهای اطلاعاتی می شود. (Cawkell, 1987) با گسترش سیستمهای پایگاهی و حجم بالای داده ها ی ذخیره شده در این سیستم ها، به ابزاری نیازاست تا بتوان این داده ها را پردازش کرد و اطلاعات حاصل از آن را در اختیار کاربران قرار داد. معمولاً کاربران پس از طرح فرضیه ای بر اساس گزارشات مشاهده شده به اثبات یا رد آن می پردازند، در حالی که امروزه به روشهایی نیازداریم که به اصطلاح به کشف دانش (Knowledge Discovery) بپردازند یعنی روشهایی که با کمترین دخالت کاربر و به صورت خ

کلمات کلیدی: داده کاوی، تولید، اشاعه و مدیریت اطلاعات و فرآیندهای اطلاعاتی

داده کاوی چیست ؟

ترجمه و تالیف: مهندس عدرا قبادی

رئیس اداره نگهداری و پشتیبانی سیستم ها

جامعه مبتنی بر اطلاعات را می توان به عنوان جامعه ای تعریف نمود که بخش غالب اجتماع به جای کارهای فیزیکی در گیرکارهای فکری هستند. در چنین جامعه ای بیشترین توجه به فعالیتهای اطلاعاتی از قبیل: فراهم آوری، پردازش، تولید، ثبت، انتقال، اشاعه و مدیریت اطلاعات مبذول می گردد و بیشترین هزینه ها صرف فرایندهای اطلاعاتی می شود. (Cawkell, 1987).

با گسترش سیستمهای پایگاهی و حجم بالای داده ها ی ذخیره شده در این سیستم ها، به ابزاری نیازاست تا بتوان این داده ها را پردازش کرد و اطلاعات حاصل از آن را در اختیار کاربران قرار داد. معمولاً کاربران پس از طرح فرضیه ای بر اساس گزارشات مشاهده شده به اثبات یا رد آن می پردازند، در حالی که امروزه به روشهایی نیازداریم که به اصطلاح به کشف دانش (Knowledge Discovery) بپردازند

یعنی روشهایی که با کمترین دخالت کاربر و به صورت خودکار الگوها و رابطه های منطقی را بیان نمایند.

یکی از روشهای بسیار مهمی که با آن می توان الگوهای مفیدی را در میان داده ها تشخیص داد، داده کاوی است، این روش که با حداقل دخالت کاربران همراه است اطلاعاتی را در اختیار آنها و تحلیل گران قرار میدهد تا براساس آنها تصمیمات مهم و حیاتی در سازمانشان اتخاذ نمایند .

باید توجه داشت که اصطلاح داده کاوی زمانی به کار برده می شود که با حجم بزرگی از داده ها ، در حد مگا یا ترا بیت ، مواجه باشیم . در تمامی منابع داده کاوی بر این مطلب تاکید شده است . هر چه حجم داده ها بیشتر و روابط میان آنها پیچیده تر باشد دسترسی به اطلاعات نهفته در میان داده ها مشکلتر می شود و نقش داده کاوی به عنوان یکی از روشهای کشف دانش ، آشکارتر می گردد.

داده کاوی از چندین رشته علمی بطور همزمان بهره میبرد نظیر : تکنولوژی پایگاه داده، هوش مصنوعی ، شبکه های عصبی، آمار، سیستم های مبتنی بر دانش، بازیابی اطلاعات و غیره . [۱] که برای پرهیز از اطاله کلام می توان آن به لحاظ تاریخی به اختصار به مراحل زیر تقسیم کرد:

مرحله اولیه: گردآوری و ایجاد پایگاه اطلاعاتی (تا دهه ۱۹۶۰)

مرحله دوم : نظامهای مدیریتی مبنی بر پایگاه اطلاعاتی (دهه ۱۹۷۰ و اوایل دهه ۱۹۸۰)

مرحله سوم : نظامهای پایگاه اطلاعاتی پیشرفته (اواسط دهه ۱۹۸۰ تا زمان حاضر)

مرحله چهارم : انبارش اطلاعات و داده کاوی (اواخر دهه ۱۹۸۰ تا به امروز)

مرحله پنجم : نظام پایگاه اطلاعاتی مبنی بر شبکه (دهه ۱۹۹۰ تا کنون)

مرحله ششم : نسل نونین نظامهای اطلاعاتی یکپارچه شده (از ۲۰۰۰ به بعد)

بدین ترتیب فعالیتی که از دهه ۱۹۶۰ شروع شده بود در دهه ۱۹۹۰ گامهای بلندی برداشت و انتظار می رود در این قرن به رشد و بالندگی خود ادامه دهد.

تعریفی از داده کاوی

بطور کلی، داده کاوی (که گاهی اوقات اکتشاف اطلاعات یا دانش نامیده میشود) عبارت از فرآیندی است که

از چشم اندازه های مختلف به تحلیل داده ها می پردازد و جمع بندی آنها را در قالب اطلاعات مفیدی ارائه میکند . این اطلاعات را میتوان برای افزایش در آمد ، کاهش هزینه ها یا هر دو به کاربرد. نرم افزار داده کاوی یکی از ابزارهای تحلیل اطلاعات است . این نرم افزار به کاربران امکان می دهد اطلاعات را از ابعاد و زوایای بسیار متفاوت تحلیل و طبقه بندی کنند و روابطی را که در آن ها شناسائی نموده اند به اجمال بیان نمایند.

به لحاظ فنی، داده کاوی عبارت از فرآیندی است که در میان حوزه های گوناگون بانکهای اطلاعاتی ارتباطی

بزرگ، همبستگی ها یا الگوهای را پیدا می کند. البته این ویژگی به معنای یکسان دانستن داده کاوی و آنالیز آماری نیست که در جدول زیر این تفاوتها آورده شده است :

داده کاوی	آنالیز آماری
به فرضیه احتیاجی ندارد.	آمارشناسان همیشه با یک فرضیه شروع به کار میکنند.
الگوریتمهای داده کاوی در ابزارها بطور اتوماتیک	آمارشناسان باید رابطه هایی را ایجاد کنند که به

روابط را ایجاد میکنند.	فرضیه آنها مربوط شود.
ابزارهای داده کاوی از انواع مختلف داده و نه فقط عددی میتوانند استفاده کنند.	آنها از داده های عددی استفاده میکنند.
داده کاوی به داده های صحیح و درست طبقه بندی شده بستگی دارد.	آنها میتوانند داده های نابجا و نادرست را در طول آنالیز تشخیص دهند.
نتایج داده کاوی آسان نیست و همچنان به متخصصان آمار برای تحلیل آنها و بیان آنها به مدیران نیاز است.	آنها میتوانند نتایج کار خود را تفسیر کنند و برای مدیران بیان کنند.

پنج ویژگی مهم داده کاوی عبارت است از :

- استخراج ، دگرگونی و بارنمودن داده های تراکنشی بر روی سیستم انبار داده ها .
- ذخیره و مدیریت داده ها در سیستم بانک اطلاعات چند بعدی.
- فراهم آوردن امکان دسترسی تحلیل گران تجاری و متخصصان تحلیل اطلاعات به داده ها .
- تحلیل داده ها با استفاده از نرم افزار کاربردی .
- معرفی نمودن ، در يك قالب بندي سودمند ، همانند گراف یا جدول

داده کاوی به چه کار می آید؟

امروزه در درجه اول شرکتها ازداده کاوی استفاده می کنند.(با توجه بسیار زیاد به مصرف کننده ، خرده فروشی، مالی ، ارتباط، وسازمانهای بازاریابی). داده کاوی این شرکتها را قادر می سازد که رابطه عوامل "درونی" (مانند قیمت ، موقع یابی فرآورده ، یا مهارت های کارمندان) ، را با عوامل "خارجی" (مانند شاخص های اقتصادی ، رقابت و آمارگیری جمعیتی مشتری) مشخص کنند؛ داده کاوی شرکت ها را قادر می سازد اثر گذاری بر مشتری ، رضایتمندی مشتری و منافع شرکت را تعیین کنند. بالاخره ، شرکتها را قادر می سازد که فشرده اطلاعات را برای دیدن داده های معاملاتی دقیق "حفاری" نمایند.

برخی از کاربردهای داده کاوی در محیطهای واقعی عبارتند از :

۱. خرده فروشی : از کاربردهای کلاسیک داده کاوی است که می توان به موارد زیر اشاره کرد :
 - تعیین الگوهای خرید مشتریان
 - تجزیه و تحلیل سبد خرید بازار
 - پیشگویی میزان خرید مشتریان از طریق فروش الکترونیکی
۲. بانکداری :
 - پیش بینی الگوهای کلاهبرداری از طریق کارتهای اعتباری

- تشخیص مشتریان ثابت
- تعیین میزان استفاده از کارتهای اعتباری بر اساس گروههای اجتماعی
- ۳. بیمه :
- تجزیه و تحلیل دعاوی
- پیشگویی میزان خرید بیمه نامه های جدید توسط مشتریان
- ۴. پزشکی :
- تعیین نوع رفتار با بیماران و پیشگویی میزان موفقیت اعمال جراحی
- تعیین میزان موفقیت روشهای درمانی در برخورد با بیماریهای صعب العلاج[۲]

نتیجه اینکه :

بسیاری از سازمانها بر معادنی از طلا تکیه زده اند. این گنجینه گرانها در شرکتهای بیمه همان داده های جمع آوری شده از بیمه گذاران، بیمه شدگان، زیاندیدگان، مقصران حادثه و انواع بیمه های فروخته شده است که می باید با بهره گیری از تکنولوژیهای جدید و ابزارهای خودکاري که بصورت هوشمند آنها را تجزیه و تحلیل می کنند، گردآوری و پردازش شده و به دانش تبدیل و به کار گرفته شوند

انبار داده ها

یکشنبه ۲۹ مهر ۱۳۸۳) تعداد دفعات خوانده شده: ۲۸۸۳)

از اواسط سال های ۱۹۸۰ نیاز به انبار داده ها به وجود آمد و دریافتند که سیستم های اطلاعاتی باید به صورت سیستم های عملیاتی و اطلاعاتی مشخص شوند. سیستم های عملیاتی از فعالیت های روزانه کسب و کار پشتیبانی می نمایند و برای پاسخگویی سریع به ارتباطات از پیش تعریف شده مناسب هستند. داده های عملیاتی ارائه بی درنگ و فعلی وضعیت کسب و کار می باشند. اما سیستم های اطلاعاتی برای مدیریت و کنترل کسب و کار به کار می روند .

این سیستم ها از تجزیه و تحلیل داده ها برای اتخاذ تصمیم درباره عملکرد آتی و آتی سازمان پشتیبانی می کنند و برای درخواست های موردی، پیچیده و به طور کلی فقط خواندنی طراحی شده اند. داده های اطلاعاتی تاریخی هستند، به عبارتی بیانگر دیدگاه ثابتی از کسب و کار در یک دوره زمانی می باشند .

ویژگی های اصلی داده های انبار داده ها

داده های موجود در انبار داده ها از سیستم های عملیاتی متنوع (نظیر سیستم های پایگاه داده ها) و منافع داده ای خارجی (نظیر پایگاه داده های آماری و WWW یکپارچه می شوند. تفاوت های ساختاری و معنایی داده ها باید پیش از یکپارچه سازی انسجام یابد. برای مثال داده ها باید مطابق با مدل داده ای یکپارچه " همگن" شوند. به علاوه، مقادیر داده ای سیستم های عملیاتی باید پاک شوند تا داده های صحیحی در انبار داده ها وارد شوند. نیاز به داده های تاریخی یکی از موارد مهم در شیوه انبار داده هاست .

داده های تاریخی برای تحلیل روند کسب و کار ضروری هستند. البته هزینه نگهداری این گونه داده ها نیز باید مورد توجه قرار گیرد. به علاوه، داده های انبار داده های ثابتی هستند، برای مثال دسترسی به DWH از نوع خواندنی است. انجام اصلاحات در این داده ها فقط هنگامی صورت می گیرد که اصلاحات داده های منبع در انبار انتشار یابند DWH. داده های دیگری به نام داده های اشتقاق یافته (derived data) دارد. این داده ها به طور صریح در منابع عملیاتی ذخیره نمی شوند، بلکه در حین بعضی از فرآیندها از داده های عملیاتی، اشتقاق می یابند. برای مثال داده های فروش را می توان در سطوح مختلف (هفتگی، ماهانه، فصلی) در انبار ذخیره نمود .

سیستم های انبار داده ها

سیستم انبار داده ها و همه مولفه هایی است که برای ساخت، دستیابی و نگهداری DWH به کار می روند. انبار داده ها بخش مرکزی سیستم انبار داده ها را تشکیل می دهد. گاهی اوقات انبار داده ها حجم عظیمی از اطلاعات را در واحد های منطقی کوچکتر به نام Data Mart نگهداری می کند. مولفه آماده سازی، مسولیت کسب یا دریافت داده ها را بر عهده دارد. این مولفه شامل همه برنامه ها و برنامه های کاربردی ای است که مسئول استخراج داده ها از منابع عملیاتی هستند. مولفه دستیابی شامل برنامه های کاربردی مختلف (OLAP) یا برنامه های کاربردی داده کاوی) است که امکان استفاده از اطلاعات ذخیره شده در انبار داده ها را فراهم می آورند .

مولفه مدیریت Metadata ، وظیفه مدیریت، تعریف و دستیابی به انواع مختلف Metadata را به عهده دارد. در اصل، Metadata "داده‌هایی درباره داده‌ها" یا "داده‌هایی است که مفهوم داده‌ها را توصیف می‌کنند". انواع مختلف Metadata در انبار داده‌ها وجود دارند. مثلاً اطلاعاتی در مورد منابع عملیاتی، ساختار داده‌های DWH و کارهایی که در حین ساخت، نگهداری و دستیابی به DWH انجام می‌شوند. نیاز به Metadata شناخته شده است. پیاده‌سازی يك DWS منسجم، کار پیچیده و دشواری است و شامل دو فاز می‌باشد. در فاز اول که پیکربندی DWS نام دارد، دیدگاه مفهومی انبار داده‌ها مطابق با نیازمندی‌های کاربر مشخص می‌شود. سپس منابع داده‌ای دخیل و روش استخراج و بارگذاری در انبار داده‌ها تعیین می‌گردد. سرانجام، درباره پایگاه داده‌های مورد نظر و روش‌های دستیابی داده‌ها تصمیم‌گیری خواهد شد. پس از بارگذاری اولیه، در فاز عملیات DWS باید داده‌های انبار داده‌ها به منظور منظم refresh شوند.

طراحی انبار داده‌ها روش‌های طراحی انبار داده‌ها امکان پردازش کارآمد query را بر روی حجم عظیمی از داده‌ها فراهم می‌آورند. نوع ویژه‌ای از الگوی پایگاه داده‌ها به نام star برای مدل‌سازی انبار داده‌های چند بعدی به کار می‌رود. در این حالت، پایگاه داده‌ها از يك جدول مرکزی واقعیت یا fact و جداول چند بعدی تشکیل شده است. جدول واقعیت حاوی tuple‌هایی است که بیانگر واقعیت‌های کسب و کار مانند فروش یا عرضه هستند. هر tuple جدول واقعیت به tuple‌های جدول چند بعدی اشاره دارد. هر tuple جدول چند بعدی نظیر محصولات، مشتریان، زمان و فروشنده را نشان می‌دهد.

انبار داده‌های مجازی

هدف انبار داده‌های مجازی، پیاده‌سازی سریع انبار داده‌ها بدون نیاز به ذخیره‌سازی و نگهداری کپی‌های متعدد از داده‌های منبع است. اغلب، انبار داده‌های مجازی به سازمانها کمک می‌کند تا به نیاز واقعی کاربران نهایی پی ببرند. کاربران نهایی می‌خواهند به طور مستقیم به داده‌های منبع بی‌درنگ با کمک ابزارهای توانمند شبکه‌ای دسترسی پیدا کنند. معایب این روش عبارتند از:

- کیفیت و سازگاری داده‌ها تضمین نمی‌شود. زیرا فعالیت‌های آماده‌سازی داده‌ها صورت نمی‌گیرند.
- به طور معمول، داده‌های تاریخی وجود ندارند.

- زمان دسترسی کاربرنهایی بسته به وجود یا عدم وجود منابع عملیاتی، بار شبکه و پیچیدگی درخواست، غیر قابل پیش‌بینی است.

معماري انباره داده از سه لايه تشكيل شده است:

-در اولين لايه اين معماری، سرويس دهنده انباره داده‌اي است که يك سيستم پایگاه داده رابطه‌اي مي باشد . اين لايه داده هاي مورد نیاز خود را از داده هاي عملياتي و منابع خارجي و فایلهاي مسطح و غيره براي ايجاد انباره داده استخراج مي کند.

-در لايه مياني يك سرويس دهنده پردازش تحلیلي برخط مي باشد که بوسيله آن مي توان مکعبهاي چند بعدي ساخت . پردازش تحلیلي برخط يك ابزار قدرتمند، سريع و مناسب براي گزارشگيري مي باشد.

-در آخرين لايه ما ابزارهاي گزارش گيري و تحليل و داده کاوي را داريم .
براي پياده سازي يك انباره داده بايد هريك از اين لايه ها به درستي پياده سازي شوند.
اخذ داده

اخذ داده از منابع مربوطه (پایگاه داده منبع) انجام مي گردد. اين مرحله بخش استخراج اطلاعات (Extract) از سري عمليات ETL است .براي انجام عمليات اخذ داده، بايد منبع اخذ داده، نحوه اخذ داده، فرمت داده هاي اخذ شده و مقاطع زماني اخذ داده ها، همچنين نحوه دسترسي به اين داده ها معلوم و مشخص باشد.

بررسي و پاکسازي داده ها

اين مرحله بخش تغيير شکل (Transform) از عمليات ETL است. بررسي و پاکسازي داده هاي استخراج شده جهت ورود به انباره داده در اين مرحله انجام مي گيرد. پس از بررسي جداول موجود، فيلدهاي موجود در جداول و محتويات فيلدهاي مذکور، کليه مشکلات داده اي در قالب ليستي ارائه مي گردند. سپس عمليات پاکسازي براي آنها انجام مي گيرد. اين عمليات غالبا در زمره يکي از موارد زیر هستند :

-حذف مقادير null

-هم مقدار سازي فيلدهاي مشابه از نظر معنا

-ايجاد فيلدهاي کمي جديد قابل بدست آمدن از روي داده هاي جدول و مورد نیاز

-يکي کردن داده ها از منابع مختلف

-خلاصه سازي سطرهاي هم معني که ايجاد افزونگي مي کنند.

-ايجاد کلید جانشين براي جداول

(Pivoting -تبدیل چند ستون به چند سطر يا بالعکس)

-تقسيم يك ستون جدول به چند ستون

طراحی انباره داده موضوعی

انباره داده موضوعی: از آنجا که کاربران مختلف با نیازهای متفاوتی وجود دارند که می‌توانند از داده‌های درون انباره داده استفاده کنند، برآوردن نیازهای تمام کاربران به وسیله یک سیستم مرکزی همیشه امکان پذیر نیست. از طرفی یک سیستم مرکزی، متمرکز بر روی داده و سیستم می‌باشد و کاربر نهایی ممکن است که بخواهد کنترل بیشتری روی محیط اطلاعاتی خود داشته باشد. راه حل این مشکلات مرکز داده ای است، که به آن انباره داده‌ای سازمانی نیز گفته می‌شود. مرکز داده ای، انباره داده خاصی است که داده‌های مورد نیاز برای یک بخش از سازمان یا کاربرهای مرتبط به آن را جمع آوری می‌کند.

طراحی مراکز داده‌ای مربوطه بنا به صلاحدید فرد خبره در قالب مدل ستاره‌ای یا دانه‌برفی و یا طرح منظومه حقایق صورت می‌گیرد. طرح ستاره‌ای: عمومی‌ترین نمونه برای مدلسازی مدل چند بعدی، طرح ستاره است. در این طرح انباره داده شامل یک جدول بزرگ مرکزی به نام جدول حقایق و یک سری جدول کوچکتر به نام جدول بعد یا جدول بعد که وابسته به جدول حقایق هستند می‌باشد.

طرح دانه برفی: این طرح، تغییر یافته طرح ستاره‌ای است بطوریکه بعضی از جداول بعد، نرمال شده‌اند. تفاوت اصلی بین طرح ستاره‌ای و طرح دانه برفی این است که جدول بعد در طرح دانه‌برفی به فرم نرمال نگهداری می‌شود تا میزان افزونگی کاهش پیدا کند. این کار باعث کاهش میزان حافظه مورد نیاز خواهد شد. البته صرفه‌جویی در فضای ذخیره‌سازی جدول بعد در مقایسه با حجم جدول حقایق ناچیز است چون تعداد اتصالاتی که برای پردازش یک گزارش باید گذارده شود در این حالت افزایش می‌یابد و مدت زمان پاسخ‌دادن به گزارش در مدل دانه برفی بیشتر از مدت زمان لازم در مدل ستاره‌ای است بنابراین غالباً طرح دانه برفی در طراحی انباره داده عمومیت طرح ستاره‌ای را ندارد مگر آنکه بنا بر صلاحدید فرد خبره بر طرح ستاره‌ای ترجیح داده شود.

طرح منظومه حقایق: هنگامی که نیاز به چندین جدول حقایق وجود دارد که دارای جداول بعدی‌های مشترک هستند، طراحی ایجاد می‌شود که به آن طرح کهکشانی یا منظومه حقایق می‌گویند. یک طرح منظومه حقایق به جداول ابعاد اجازه می‌دهد که بین جداول حقایق مشترک باشند.

وارد سازي داده‌هاي پاكسازي شده به انباره داده

با توجه به فرمت داده‌هاي اخذ شده، وارد سازي داده‌هاي پاكسازي شده به انباره داده با اجرائي اسكربت مربوطه - بخش Load از عمليات ETL.

تحليل داده با ابزار پردازش تحليلي برخط و داده‌كاوي

عمليات پردازش تحليلي برخط بر انباره داده‌هاي موضوعي اعمال مي‌گردد. بسياري فكر مي‌كنند كه داده كاوي و OLAP دو چيز مشابه هستند در اين بخش سعي مي‌كنيم اين مسئله را بررسي كنيم و همانطور كه خواهيم ديد اين دو ابزار هاي كاملا متفاوت مي‌باشند كه مي‌توانند همديگر را تكميل كنند.

OLAP جزئي از تكنيك‌هاي تصميم گيري مي‌باشد. سيستم هاي سنتي گزارش گيري و پايگاه داده اي آنچه را كه در پايگاه داده بود توضيح مي‌دادند حال آنكه در OLAP هدف بررسي دليل صحت يك فرضيه است. [Error! Reference source not found.]

بدين معني كه كاربر فرضيه اي در مورد داده ها و روابط بين آنها ارائه مي‌كند و سپس به وسيله ابزار OLAP با انجام چند Query صحت آن فرضيه را بررسي مي‌كند.

اما اين روش براي هنگامي كه داده ها بسيار حجيم بوده و تعداد پارامترها زياد باشد نميتواند مفيد باشد چون حدس روابط بين داده ها كار سخت و بررسي صحت آن بسيار زمانبر خواهد بود. تفاوت داده كاوي با OLAP در اين است كه داده كاوي برخلاف OLAP براي بررسي صحت يك الگوي فرضي استفاده نمي‌شود بلكه خود سعي مي‌كند اين الگوها را كشف كند.

درنتيجه داده كاوي و OLAP مي‌توانند همديگر را تكميل كنند و تحليل گر مي‌تواند به وسيله ابزار OLAP يك سري اطلاعات كسب كند كه در مرحله داده كاوي مي‌تواند مفيد باشد و همچنين الگوها و روابط كشف شده در مرحله داده كاوي مي‌تواند درست نباشد كه با اعمال تغييرات در آنها مي‌توان به وسيله OLAP بيشتر بررسي شوند.

با گسترش شگرف اينترنت و استفاده روزافزون از آن در جهت ارايه و يا كسب اطلاعات، شاهد حجم انبوهي از اسناد و مقالات بر- خط هستيم كه بعنوان يكي از مشخصات بارز زندگي مدرن امروزي، تحت عنوان افزونگي اطلاعاتي مطرح مي‌گردد. در اين ميان دسترسي سريع و صحيح به منابع مهم و مورد علاقه، يكي از دغدغه هاي استفاده كنندگان از اين منبع اطلاعاتي بسيار بزرگ است. آنچه امروزه از اهميت بسيار زيادي برخوردار گرديده ، كمبود يا نبود اطلاعات نيست بلكه كمبود روشهايي در جهت يافت و بهره برداري از اطلاعات در دسترس به نحوي بهينه است. بعنوان مسئله اي آرمانی تر به دنبال روشهايي هستيم تا از اطلاعات موجود به كسب دانش پرداخته،

احتمالاً به ارایه مسایل جدیدی بپردازد که قبل از آن مشخص نبوده است. متن کاوی، بعنوان روشی در استخراج دانش از متون، یکی از موضوعات مهم در گستره ای از اعمال مدیریت اطلاعات است. در این میان آنچه از اهمیت فوق العاده ای برخوردار است ارایه راه کارهایی برای مواجهه با این حجم عظیم اطلاعاتی و استفاده بهینه از اطلاعات در جهت خلق دانش، تولید سینرجی و در نهایت افزایش خرد جمعی است. در سالهای اخیر اهمیت متون به عنوان منابع با پتانسیل اطلاعاتی بسیار بالا به نحو گسترده ای مورد توجه قرار گرفته به طوری که کشف دانش از متون به عنوان یکی از مهمترین فعالیتهای محققین حوزه هوش مصنوعی و فناوری اطلاعات قرار گرفته است. تحقیقات بسیاری صورت گرفته اما محدوده فعالیت بقدری گسترده است که نیازمند توجه بیشتری می باشد. امروزه محققان به این مسئله معترفند که با وجود انجام تحقیقات بی وقفه در زمینه کاری خود، نمی توانند همزمان با پیشرفت دانش، معلومات خود را به روز نگاه دارند. بعنوان مثال بانک اطلاعاتی Medline در حال حاضر حاوی ۱۰ میلیون چکیده مقاله است و هر هفته بین هفت تا هشت هزار چکیده مقاله به این بانک اطلاعاتی افزوده می شود. در این بین شاید همه مقالات مربوط به يك دانش خاص نباشند، اما تعداد مقالات تخصصی که در حوزه تحقیق يك دانش خاص قرار می گیرد به اندازه ای است که يك نفر نمی تواند ادعا کند همه آنها را مطالعه کرده است بعلاوه نقش مطالعات عمیق و گسترده و استخراج ایده ها و دانش جدید از مطالب مطالعه شده بر کسی پوشیده نیست. در این میان اینترنت بعنوان بزرگترین منبع اطلاعاتی همگانی، تشکیل یافته از صد ها میلیون صفحه اطلاعات است که به جهت همگانی بودن آن و نبود آینده نگری کافی در زمان تشکیل و رشد آن، متحمل نگاهداری اطلاعات نویسندگان، محققان، اندیشمندان و غیره به همان نحوی که آنها می نوشتند گردید. نبود يك استاندارد همه جانبه و دقیق در تنظیم متون و قرار گیری این مجموعه عظیم بصورتی غیر ساختیافته و یا بعضاً نیمه ساختیافته، جامعه اطلاعاتی را دچار نوعی سردرگمی و مشکل در دستیابی به اطلاعات مورد نیاز کرده بطوری که برای یافتن مطالب مورد نظر خود متحمل هزینه های زمانی بسیاری می گردند. محققان به ارایه راه کارهایی برای ساخت یافته کردن اطلاعات نمودند و با ارایه زبانهای نشانه گذاری استاندارد نظیر XML تا حد زیادی جلوی این از هم پاشیدگی اطلاعاتی را گرفتند اما آنچه همچنان باقی است وجود بسیاری از متون غیر ساخت یافته می باشد؛ در همین راستا ارایه ابزارهایی که با بررسی متون بتوانند تحلیلی روی آنها انجام دهند منجر به شکل گیری زمینه ای جدید در هوش مصنوعی و فناوری اطلاعات گردیده که به یادگیری متن معروف است.

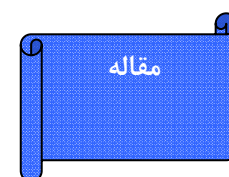
این حوزه تمام فعالیتهایی که به نوعی به دنبال کسب دانش از متن هستند را شامل می گردد. آنالیز داده های متنی توسط تکنیکهای یادگیری ماشین، بازیابی اطلاعات هوشمند، پردازش زبان طبیعی یا روشهای مرتبط دیگر همگی در زمره مقوله یادگیری متن قرار می گیرند. یکی از روشهایی که ذکر گردید، استفاده از تکنیکهای یادگیری ماشین در زمینه پردازش متن است، مسئله قابل تامل این است که این تکنیکها در ابتدا در مورد داده های ساخت یافته به کار گرفته شدند و علمی به نام داده کاوی را بوجود آوردند. داده های ساخت یافته به داده هایی اطلاق می گردد که بطور کاملاً مستقل از همدیگر ولی یکسان از لحاظ ساختاری در يك محل گردآوری شده اند. انواع بانکهای اطلاعاتی را می توان نمونه هایی از این دسته اطلاعات نام برد. در اینصورت مسئله داده کاوی عبارت از کسب اطلاعات و دانش از این مجموعه ساخت یافته، اما در مورد متون که عمدتاً غیر ساخت یافته یا نیمه ساخت یافته هستند ابتدا باید توسط روشهایی، آنها را ساختارمند نمود و

سپس از این روشها برای استخراج اطلاعات و دانش از آنها استفاده کرد. به هر حال استفاده از داده کاوی در مورد متن خود شاخه ای دیگر را در علوم هوش مصنوعی بوجود آورد به نام متن کاوی . از جمله فعالیتهای بسیار مهم در این زمینه، طبقه بندی (دسته بندی) متن می باشد. طبقه بندی متن، یعنی انتساب اسناد متنی بر اساس محتوی به یک یا چند طبقه از قبل تعیین شده، یکی از مهمترین مسایل در متن کاوی است؛ مرتب کردن بلادرنگ نامه های الکترونیکی یا فایلها در سلسله مراتبی از پوشه ها، تشخیص موضوع متن، جستجوی ساختیافته و/ یا پیدا کردن اسنادی که در راستای علایق کاربر میباشد، از جمله کاربردهای میحث طبقه بندی (دسته بندی- کلاس بندی) متن است. در بسیاری از موارد، افراد حرفه ای آموزش دیده، برای طبقه بندی متون جدید به کار گرفته می شوند. این فرآیند بسیار زمان بر و پرهزینه است و لذا کاربرد خود را محدود می سازد، به همین منظور علاقه روزافزونی به توسعه فناوری هایی در دسته بندی خودکار متن ابراز میشود.

در هر حال در جوامع اطلاعاتی امروزی آنچه از اهمیت روزافزونی برخوردار است، اطلاعات و تبادل آن است و در این راستا به توسعه فناوری های مرتبط پرداخته می شود، اما یک مرحله کاملاً جدید تر و کاملاً مورد توجه جوامع فرا صنعتی، خلق دانش جدید از اطلاعات قبلی است که این جوامع آنرا کلید موفقیت خود در آینده دانسته و به سختی در این زمینه فعالیت می نمایند. بر ما است تا ضمن ارتقای فناوری اطلاعات در کشور و ایجاد زیر ساختهای لازمه در اسرع وقت، به اینگونه مسائل جدی تر که در زمره Information High Technology قرار می گیرند، بپردازیم.



شماره ۱۸۵- اسفند ۸۴



داده کاوی، مفهوم و کاربرد آن در آموزش عالی

احمد سعیدی

دانشجوی دکتری اقتصاد و مدیریت مالی آموزش عالی

مقدمه

از هنگامی که رایانه در تحلیل و ذخیره سازی داده ها بکار رفت (۱۹۵۰) پس از حدود ۲۰ سال، حجم داده ها در پایگاه داده ها دو برابر شد. ولی پس از گذشت دو دهه و همزمان با پیشرفت فن آوری اطلاعات (IT) هر دو سال یکبار حجم داده ها، دو برابر شد. همچنین تعداد پایگاه داده ها با سرعت بیشتری رشد نمود. این در حالی است که تعداد متخصصین تحلیل داده ها و آمارشناسان با این سرعت رشد نکرد. حتی اگر چنین امری اتفاق می افتاد، بسیاری از پایگاه داده ها چنان گسترش یافته اند که شامل چندصد میلیون یا چندصد میلیارد رکورد ثبت شده هستند و امکان تحلیل و استخراج اطلاعات با روش های معمول آماری از دل انبوه داده ها مستلزم چند روز کار با رایانه - های موجود است. حال با وجود سیستم های یکپارچه اطلاعاتی، سیستم های یکپارچه بانکی و تجارت الکترونیک، لحظه به لحظه به حجم داده ها در پایگاه داده های مربوط اضافه شده و باعث به وجود آمدن انبارهای (توده های) عظیمی از داده ها شده است به طوری که ضرورت کشف و استخراج سریع و دقیق دانش از این پایگاه داده ها را بیش از پیش نمایان کرده است (چنان که در عصر حاضر گفته می شود «اطلاعات طلاست»).

هم اکنون در هر کشور، سازمان ها، شرکت ها و ... برای امور بازرگانی، پرسنلی، آموزشی، آماری و ... پایگاه داده ها ایجاد یا خریداری شده است، به طوری که این پایگاه داده ها برای مدیران، برنامه ریزان، پژوهشگران و ... جهت تصمیم گیری های راهبردی، تهیه گزارش های مختلف، توصیف وضعیت جاری خود و ... می تواند مفید باشد. داده کاوی یا استخراج و کشف سریع و دقیق اطلاعات با ارزش و پنهان از این پایگاه داده ها از جمله اموری است که هر کشور، سازمان و شرکتی به منظور توسعه علمی، فنی و اقتصادی خود به آن نیاز دارد.

در کشور ما نیز سازمان ها، شرکت ها و مؤسسات دولتی و خصوصی به طور فزاینده ولی آهسته در حال ایجاد یا خرید نرم افزارهای پایگاه داده ها و مکانیزه کردن سیستم های اطلاعات خود هستند، همچنین با توجه به فصول دهم و یازدهم قانون برنامه سوم توسعه در خصوص داد و ستدهای الکترونیکی و همچنین تأکید بر برخورداری

کشور از فن آوری های جدید اطلاعات برای دستیابی آسان به اطلاعات داخلی و خارجی، دولت مکلف شده است امکانات لازم برای دستیابی آسان به اطلاعات، زمینه سازی برای اتصال کشور به شبکه های جهانی و ایجاد زیر ساخت های ارتباطی و شاهراه های اطلاعاتی فراهم کند. واضح است این امر باعث ایجاد پایگاه های عظیم داده ها شده و ضرورت استفاده از داده کاوی را بیش از پیش نمایان می سازد.

سابقه داده کاوی

داده کاوی و کشف دانش در پایگاه داده ها از جمله موضوع هایی هستند که همزمان با ایجاد و استفاده از پایگاه داده ها در اوایل دهه ۸۰ برای جستجوی دانش در داده ها شکل گرفت.

شاید بتوان لوول (۱۹۸۳) را اولین شخصی دانست که گزارشی در مورد داده کاوی تحت عنوان « شبیه سازی فعالیت داده کاوی » ارائه نمود. همزمان با او پژوهشگران و متخصصان علوم رایانه، آمار، هوش مصنوعی، یادگیری ماشین و ... نیز به پژوهش در این زمینه و زمینه های مرتبط با آن پرداخته اند.

پژوهش جدی روی موضوع داده کاوی از اوایل دهه ۹۰ شروع شد. پژوهش ها و مطالعه های زیادی در این زمینه صورت گرفته، همچنین سمینارها، دوره های آموزشی و کنفرانس هایی نیز برگزار شده است. نتایج پایه های نظری داده کاوی در تعدادی از مقاله های پژوهشی آورده شده است. مثلاً سال ۱۹۹۱ پیاتتسکی و شاپیرو ۲ « استقلال آماری قاعده ها در داده کاوی » را بررسی نموده اند. سال ۱۹۹۵ هافمن و نش استفاده از داده کاوی و داده انبار ۳ توسط بانک های آمریکا را بررسی نموده و بیان کردند که چگونه این سیستم ها برای بانک های آمریکا قدرت رقابت بیشتری ایجاد می کنند. چت فیلد مشکلات ایجاد شده توسط داده کاوی را بررسی نمود و همچنین مقاله ای تحت عنوان « مدل های خطی غیر دقیق داده کاوی و استنباط آماری » ارائه نمود. هندری نیز دیدگاه اقتصاد سنجی روی داده کاوی را تهیه کرد. در این سال انجمن داده کاوی همزمان با اولین کنفرانس بین المللی « کشف دانش و داده کاوی » شروع به کار کرد. این کنفرانس

توسعه یافته چهار دوره آموزشی بین المللی در پایگاه های داده در سال ۱۹۸۹ تا ۱۹۹۴ بود. انجمن مذکور، یک سازمان علمی به نام ACM- SIGKDD را ایجاد نمود. سال ۱۹۹۶ ایمیلنسکی ۴ و منیلا ۵ دیدگاهی از داده کاوی به عنوان «پرس و جو کننده از پایگاه های استنتاجی ۶» را پیشنهاد کردند. فایاد، پیاتتسکی – شاپیرو، اودوراسامی پیشرفت های کشف دانش و داده کاوی را عنوان کردند. در سال ۱۹۹۷ منیلا خلاصه ای از مطالعه روی اساس داده کاوی ارائه نمود. باربارا و همکاران نیز دیدگاه کاهش داده ها روی داده کاوی را در گزارش کاهش داده های نیوجرسی ارائه نمودند. همچنین می توان برای کاربرد داده کاوی در مدیریت مالی می توان، تحلیل داده های مالی و مدل سازی مالی بنینگا و چاچ کز و هیگینز ۷ را ملاحظه کرد فریدمن نیز مقاله ای در ارتباط با مفهوم آمار و داده کاوی ارائه نمود. سال ۱۹۹۸ هند ۸ مقاله ای تحت عنوان «داده کاوی : آمار یا بیشتر؟» ارائه نمود. کلینبرگ ۹ پائودیمیتریو و راغان ۱۰ دیدگاه اقتصاد سنجی روی داده کاوی و عملکرد داده کاوی به عنوان یک مسئله بهینه را ارائه نمودند. در این سال نیز کنفرانس های ناحیه ای و بین المللی در مورد داده کاوی برگزار شد که از جمله می توان به کنفرانس آسیا و اقیانوسیه درباره کشف دانش و داده کاوی اشاره کرد. سال ۲۰۰۰ هند و همکاران و اسمیت بحث های مقایسه ای بین آمار و داده کاوی را ارائه کردند. سری و استاوا، کولی، رش پاند و تن استفاده از وب در کاوش داده ها و کاربردهای آن را ارائه کردند. سال ۲۰۰۲ کلادیو کانورسانو و همکاران «مدل آمیخته چندگانه جمع پذیر تعمیم یافته» برای داده کاوی را بررسی نمودند. پائلو و گیانلوکاپاسرون، «داده کاوی ساختارهای پیوند برای مدل رفتار مصرف کننده» را ارائه نمودند.

-
- 4- Imielnski
 - 5 - Manila
 - 6 - Inductive databases
 - 7 - Benninga, Czaczkes, Higgins
 - 8 - Hand
 - 9 - Kleinberg
 - 10 - Paodimitriou , Raghavan

مفهوم داده کاوی

عبارت داده کاوی مترادف با یکی از عبارات های استخراج دانش، برداشت اطلاعات، واری داده ها و حتی لایروبی کردن داده هاست که در حقیقت کشف دانش در پایگاه داده ها (KDD) را توصیف می کند. بنابراین ایده ای که مبنای داده کاوی است یک فرآیند با اهمیت از شناخت الگوهای بالقوه مفید، تازه و درنهایت قابل درک در داده هاست. واژه کشف دانش در پایگاه داده ها در اوایل دهه ۸۰ در مراجعه به مفهوم کلی، گسترده، سطح بالا و به دنبال جستجوی دانش در اطلاعات شکل گرفته است. داده کاوی کاربرد سطح بالای فنون و ابزار بکار برده شده برای معرفی و تحلیل داده های تصمیم گیرندگان است. اصطلاح داده کاوی را آمار شناسان، تحلیل گران داده ها و انجمن سیستم های اطلاعات مدیریت به کار برده اند در حالی که پژوهشگران یادگیری ماشین و هوش مصنوعی از KDD بیشتر استفاده می کنند. در ادامه چند تعریف از داده کاوی ارائه می شود.

۱- «داده کاوی یا به تعبیر دیگر کشف دانش در پایگاه داده

ها، استخراج غیر بدیهی اطلاعات بالقوه مفید از روی داده هایی است که قبلاً ناشناخته مانده اند. این مطلب برخی از روش های فنی مانند خوشه بندی، خلاصه سازی داده ها، فراگیری قاعده های رده بندی، یافتن ارتباط شبکه ها، تحلیل تغییرات و کشف بی قاعدگی را شامل می شود» (پیاتسکی شاپیرو، مائوس کریستوفر)

۲- «داده کاوی در حقیقت کشف ساختارهای جالب توجه،

غیر منتظره و با ارزش از داخل مجموعه وسیعی از داده ها می باشد و فعالیتی است که اساساً با آمار و تحلیل دقیق داده ها منطبق است» هند (۱۹۹۸)

۳- «داده کاوی فرآیند کشف رابطه ها، الگوها و روندهای

جدید معنی داری است که به بررسی حجم وسیعی از اطلاعات ذخیره

شده در انبارهای داده با فناوری های تشخیص الگو (مانند ریاضی و آمار)

می پردازد». (سایت ۱۲ <http://www.spss.com>)

کشف دانش در پایگاه داده ها در جهت کشف اطلاعات مفید از مجموعه بزرگ داده هاست. دانش کشف شده می تواند قاعده ای باشد تا ویژگی های داده ها، الگوهایی که به طور متناسب رخ می دهند، خوشه بندی موضوع های درون پایگاه داده ها و غیره را توصیف می کند.

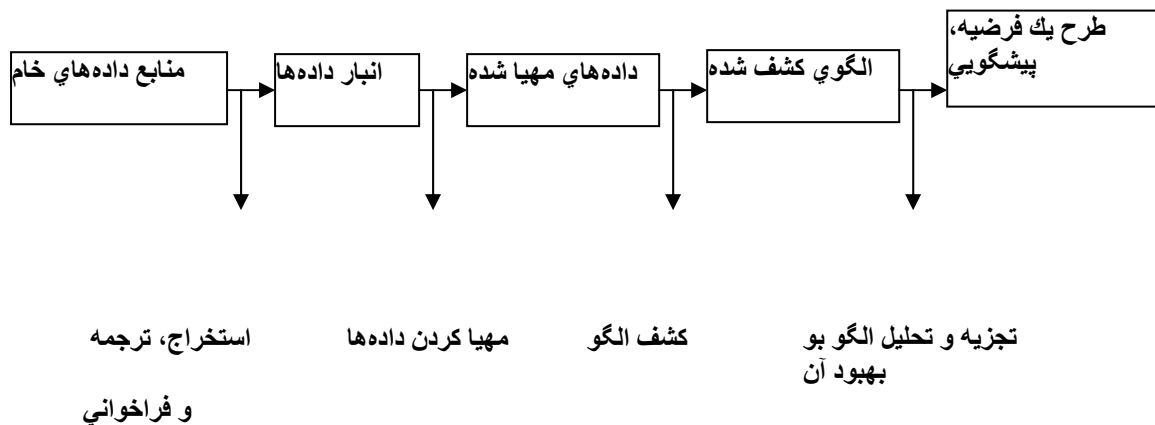
یک کاربر سیستم **KDD** بایستی درک بالایی از قلمرو داده ها به منظور انتخاب زیر مجموعه صحیحی از داده ها، رده مناسبی از الگوها و معیار خوبی برای الگوهای جالب داشته باشد. بنابراین سیستم **KDD** باید ابزارهایی با اثر تعاملی داشته باشد نه سیستم های تجزیه و تحلیل خودکار. لذا کشف دانش از پایگاه داده ها باید مثل یک فرآیند شامل گام های زیر باشد:

- ۱- درک قلمرو
- ۲- آماده کردن مجموعه داده ها
- ۳- کشف الگوها (داده کاوی)
- ۴- پردازش بعد از کشف الگو
- ۵- استفاده از نتایج .

فرآیند داده کاوی

می توان فرآیند داده کاوی را طی مراحل زیر به صورت نمودار نشان داد.

۱- این سایت یکی از معتبرترین سایت ها در زمینه آمار و داده کاوی است.



در فرآیند بالا، داده‌های خام از منابع مختلفی جمع‌آوری می‌شوند و از طریق استخراج، ترجمه و فرآیندهای بازخوانی به انبار داده‌ها وارد می‌شوند. در بخش مهیاسازی داده‌ها، داده‌ها از انبار خارج شده و به صورت یک فرمت مناسب برای داده‌کاوی درمی‌آیند. در بخش کشف الگو با روش‌های داده‌کاوی برای پاسخ به سؤال‌های خاصی که به ذهن می‌رسند، الگوریتم‌هایی را استخراج می‌کنند و از این الگوریتم‌ها برای ساخت الگو استفاده می‌شود. در بخش تجزیه و تحلیل الگو، 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.

نرم‌افزارهای داده‌کاوی

طی سال‌های گذشته جریان سریعی از تمایل به داده‌کاوی در بازارهای نرم‌افزاری به وجود آمده است. بیشتر کاربران نرم‌افزارهای داده‌کاو با تفکر استفاده تجاری از این نرم‌افزارها، خواهان استفاده از آن شده‌اند. نرم‌افزارهای داده‌کاو معمولاً سه روش مختلف را برای استفاده از داده‌کاوی به کار می‌برند. (۱) اکتشاف (۲) استفاده از مدل‌های پیشگویی (۳) استفاده از آنالیز بحث و جدل.

اکتشاف، فرآیند جستجو در داده‌هاست تا الگوهای مخفی موجود در داده‌ها را بدون هیچ ایده‌ای از پیش تعیین شده‌ای مشخص نماید. در نرم‌افزارهای داده‌کاوی مبتنی بر

مدل‌های پیشگویی، الگوهایی که از یک بانک داده کشف می‌شوند، برای پیش‌بینی آینده به کار می‌روند. مدل‌های پیش‌بینی به کاربر اجازه می‌دهند تا داده‌های نامشخص را به کار ببرد و این مقادیر نامشخص توسط نرم‌افزار کشف شود.

در مدل‌های جدلی نیز الگوهای یافت شده از داده‌ها برای تعیین مقادیر غیرعادی به کار می‌رود. برای تعیین مقادیر غیر عادی، ابتدا می‌بایست مقادیر عادی شناخته شود تا بر این اساس مقادیر غیرعادی و منحرف شناخته شوند.

نرم‌افزارهای داده‌کاوی در حال حاضر از فعالیت کمتری نسبت به سایر نرم‌افزارهای هوشمند برخوردار هستند. با این وجود فعالیت تجاری این نرم‌افزار را می‌توان در شش بخش کلی، دسته‌بندی داده‌ها، برآورد مقادیر نامشخص، پیش‌بینی مقادیر نامشخص، گروه‌بندی تقریبی داده‌ها، خوشه‌بندی داده‌ها و تشریح روابط بین داده‌ها تقسیم کرد. داده‌کاوی و مدیریت دانش

اگر چه دانش به طور انحصاری محصول فناوری اطلاعات نیست، ولی فناوری اطلاعات به طور لاینفکی در ایجاد دانش و فرآیند مدیریت دانش از سال‌های اول مشارکت داشته است. امروزه مدیریت دانش از مسئولیت‌های فناوری اطلاعات به شمار می‌رود. زیرا در جمع‌آوری، تبدیل دانش و انتقال داده‌ها، اطلاعات و دانش نقش کلیدی دارد.

از منظر مدیریت دانش، هدف داده‌کاوی، کشف دانش سازمانی پنهان در اطلاعات خام است. اینگونه نیست که هر بینش حاصل از داده‌کاوی دانش می‌سازد، بلکه در عوض بسیاری از نتایج به دست آمده، اطلاعات مدیریت، یا هوش سازمانی است. مثلاً در سازمان‌های تجاری، دانش با ارزش

مورد مشتری، محصول و بازار را می‌توان از طریق داده‌کاوی به دست آورد. داده‌کاوی ابزار مفیدی برای مدیران دانش است که کشف را با تحلیل تلفیق می‌کنند. تلفیقی که اغلب منجر به ایجاد دانش می‌شود.

کاربرد داده‌کاوی در آموزش عالی

با توجه به اینکه آموزش عالی همواره با داده‌ها و اطلاعات بسیار زیادی در مورد دانشگاه‌ها، دانشجویان، اعضای هیئت علمی، پرسنل، منابع مادی و... روبروست و در اکثر

مواقع این داده‌ها می‌تواند حامل اطلاعات و الگوهای باارزشی باشند، لذا به نظر می‌رسد یکی از مهمترین کاربردهای داده‌کاوی در آموزش عالی است. امروز بانک‌های اطلاعاتی وسیعی از ویژگی‌های دانشجویان موجود است که اطلاعات مربوط به ویژگی‌های خانوادگی، تحصیلی و ... را شامل می‌شود. پیدا کردن الگوها و دانش نهفته در این اطلاعات می‌تواند به تصمیم‌گیرندگان عرصه آموزش عالی کمک شایانی بکند. استفاده از تکنیک‌های پیشرفته داده‌کاوی مانند خوشه‌بندی، طبقه‌بندی، و ... می‌تواند در طبقه‌بندی دانشگاه‌ها، یافتن الگوهای خاص و با ارزش در مورد دانشجویان موفق، یافتن یک برنامه یا روش موفق تدریس، یافتن نقاط بحرانی در مدیریت مالی دانشگاه‌ها و موارد دیگر کاربرد داشته باشد.

نتیجه‌گیری

شرکت‌ها، سازمان‌ها، دانشگاه‌ها و مؤسسات آموزش عالی امروزی غرق در انبوه داده‌ها و اطلاعاتی هستند که استفاده از آنها در بیشتر موارد محدود به انجام کارهای جاری می‌باشد و هنوز از داده‌ها در تصمیم‌گیری استراتژیک استفاده نمی‌شود. داده‌کاوی که استفاده از آن روز به روز توسعه می‌یابد می‌تواند به استفاده از اطلاعات موجود در مؤسسات و مراکز آموزش عالی در زمینه‌های تصمیم‌گیری استراتژیک منجر شود.

منابع:

- ۱- مهریزی، حائری، علی اصغر، «داده‌کاوی: مفاهیم، روش‌ها و کاربردها» (۱۳۸۲) پایان‌نامه کارشناسی ارشد آمار اقتصادی و اجتماعی، دانشکده اقتصاد، دانشگاه علامه طباطبائی.
- ۲- زعفریان، رضا و زعفریان، قاسم، «مروری بر داده‌کاوی» (۱۳۸۰) فصلنامه صنایع، شماره ۲۹
- ۳- شاه‌سمندی، پرستو «داده‌کاوی در مدیریت ارتباط با مشتری» (۱۳۸۴)، مجله تدبیر شماره ۱۵۶.

۴- گودرزی، حمیدرضا، مترجم «داده‌کاوی چیست»، نشریه گزیده مطالب
آماري، مرکز آمار ایران، شماره ۵۲.

5) Hand. D.J (1998): "Review of Data mining", The American
statistician, 52, 112-118.