

داده کای در مدیریت ارتباط با مشتری

هومن حری

دانشجوی مهندسی فن آوری اطلاعات

horrihooman@gmail.com

چکیده

امروزه با گسترش سیستم های پایگامی و حجم بالای داده های ذخیره شده در این سیستم ها، نیاز به ابزاری است تا بتوان داده های ذخیره شده پردازش کرد و اطلاعات حاصل از این پردازش را در اختیار کاربران قرار داد.

با استفاده از پرسش های ساده در SQL و ابزارهای گوناگون گزارش گیری معمولی، می توان اطلاعاتی را در اختیار کاربران قرار داد تا بتوانند به نتیجه گیری در مورد داده ها و روابط منطقی میان آنها بپردازند اما وقتی که حجم داده ها بالا باشد، کاربران هر چقدر حرفه ای و با تجربه باشند نمی توانند الگوهای مفید را در میان حجم انبوه داده ها تشخیص دهند و یا اگر قادر به این کار هم باشند، هزینه عملیات از نظر نیروی انسانی و مالی بسیار بالا است.

بنابراین می شود گفت که در حال حاضر یک تغییر الگو از مدل سازی و تحلیل های کلاسیک برپایه اصول اولیه به مدل های در حال پیشرفت و تحلیل های مربوط بطور مستقیم از داده ها وجود دارد.

داده کاوی یکی از مهمترین این روشها است که به وسیله آن الگوهای مفید در داده ها با حداقل دخالت کاربران شناخته می شوند و اطلاعاتی را در اختیار کاربران و تحلیل گران قرار می دهند تا براساس آنها تصمیمات مهم و حیاتی در سازمانها اتخاذ شوند.

در داده کاوی از بخشی از به نام تحلیل اکتشافی داده ها استفاده می شود که در آن بر کشف اطلاعات نهفته و ناشناخته از درون حجم انبوه داده ها تاکید می شود بنابراین می توان گفت در داده کاوی تئوریهای پایگاه داده ها، هوش مصنوعی، یادگیری ماشین و علم آمار را در هم می آمیزند تا زمینه کاربردی فراهم شود.

باید توجه داشت که اصطلاح داده کاوی زمانی به کار برده می شود که با حجم بزرگی از داده ها در حد گیگابایت یا ترابایت، مواجه باشیم که از این نظر یکی از بزرگترین بازارهای هدف، انبارجامع داده ها، مراکز داده و سیستم های پشتیبانی تصمیم برای بدست آوردن تخصص هایی در صنایعی مثل شبکه های توزیع مویرگی، تولیدف مخابرات، بیمه و ... می باشد.

مدیریت ارتباط با مشتری به همه فرآیندها و فناوریهایی گفته می شود که در شرکتها و سازمانها برای شناسایی، ترغیب، گسترش، حفظ و ارائه خدمت به مشتریان به کار می رود.

سازمانها با استفاده از CRM می توانند چرخه فروش را کوتاهتر و وفاداری مشتری به ایجاد روابط نزدیکتر و درآمد را افزایش دهند. سیستم مدیریت روابط با مشتری می تواند کمک کند تا مشتریان موجود حفظ شوند و مشتریان جدیدی جذب شوند. سازمانها برخی روشهایی را شامل مدیریت ارتباط با مشتری،

تحلیل ارزش مشتری، استراتژی سازمانی و ساز و کارهای خدماتی که کارایی ارتباطات مشتری را بهبود می‌دهد بکار می‌برند. مدیریت ارتباط با مشتری استراتژی‌ای برای کسب مشتریان جدید و نگهداشتن آنها است. مدیریت ارتباط با مشتری عملیاتی شامل تمام فعالیتهای مرتبط با مشتریان بی واسطه همچون شرکتها می‌باشد.

ارتباط بین این دو (داده کاوی و مدیریت ارتباط با مشتری) خیلی مهم و اساسی می‌باشد زیرا شرکت‌ها با داده کاوی خواسته‌های مشتریان خود را می‌شناسند و بر اساس با آن به مشتریان خود سرویس می‌دهند در این مقاله در بخش اول با مقدمه‌ای از داده کاوی و تاریخچه‌ای در این مورد را بررسی می‌کنیم و در بخش بعد با داده کاوی آشنا می‌شویم بعد با وب کاوی که اساس بسیاری از شرکت‌ها برای استخراج داده‌های خام از پرتال‌های خود هستند آشنا می‌شویم در بخش چهارم با مدیریت ارتباط با مشتری آشنا می‌شویم و در فصل آخر ارتباط بین این دو را بررسی می‌کنیم

فصل اول

مقدمه

۱-۱ مقدمه :

از هنگامی که رایانه در تحلیل و ذخیره سازی داده ها بکار رفت (1950) پس از حدود 20 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 شده و همچنین تعداد پایگاه داده ها با سرعت بیشتری رشد نمود. این در حالی است که تعداد متخصصین تحلیل داده ها 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. 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 عظیمی از داده ها شده است.

این واقعیت، 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، شرکت و غیره برای امور بازرگانی، پرسنلی، آموزشی، آماری و غیره پایگاه داده ها ایجاد یا خریداری شده است. به طوری که این پایگاه داده ها برای مدیران، برنامه ریزان، پژوهشگران جهت، تصمیم گیری های راهبردی، تهیه گزارش های مختلف، توصیف وضعیت جاری خود و سایر اهداف می تواند مفید باشد. بسیاری از این داده ها از نرم افزارهای تجاری، مثل کاربردهای مالی، ERPها، CRMها و web log ها، می آیند. نتیجه این جمع آوری داده ها این می شود که در سازمانها، داده ها غنی ولی دانش ضعیف، است. جمع آوری داده ها، بسیار انبوه می شود و بسرعت اندازه آن افزایش می یابد و استفاده عملی از داده ها را محدود می سازد.

داده کاوی استخراج و تحلیل مقدار زیادی داده بمنظور کشف قوانین و الگوهای معنی دار در آنهاست. هدف اصلی داده کاوی، استخراج الگوهایی از داده ها، افزایش ارزش اصلی آنها و انتقال داده ها بصورت دانش است.

داده کاوی، بهمراه OLAP، گزارشگری تشکیلات اقتصادی (Enterprise reporting) و ETL، یک عضو کلیدی در خانواده محصول Business Intelligence (BI) است. حوزه های مختلفی وجود

دارد که در آنها حجم بسیاری از داده در پایگاه داده‌های متمرکز یا توزیع شده ذخیره می‌شود.

برخی از آنها به قرار زیر هستند: **Error! Reference source not found.**

- 1- کتابخانه دیجیتال: یک مجموعه سازماندهی شده از اطلاعات دیجیتال که بصورت متن در پایگاه داده‌های بزرگی ذخیره می‌شوند.
- 2- آرشیو تصویر: شامل پایگاه داده بزرگی از تصاویر به شکل خام یا فشرده.
- 3- اطلاعات زیستی: بدن هر انسانی از 50 تا 100 هزار نوع ژن یا پروتئین مختلف ساخته شده است. اطلاعات زیستی شامل تحلیل و تفسیر این حجم عظیم داده ذخیره شده در پایگاه داده بزرگی از ژنهاست.
- 4- تصاویر پزشکی: روزانه حجم وسیعی از داده‌های پزشکی به شکل تصاویر دیجیتال تولید می‌شوند، مانند EKG، MRI، ACT، SCAN و غیره. اینها در پایگاه داده‌های بزرگی در سیستم‌های مدیریت پزشکی ذخیره می‌شوند.
- 5- مراقبت‌های پزشکی: بجز اطلاعات بالا، یکسری اطلاعات پزشکی دیگری نیز روزانه ذخیره می‌شود مانند سوابق پزشکی بیماران، اطلاعات بیمه درمانی، اطلاعات بیماران خاص و غیره.
- 6- اطلاعات مالی و سرمایه‌گذاری: این اطلاعات دامنه بزرگی از داده‌ها هستند که برای داده‌کاوی بسیار مطلوب می‌باشند. از این قبیل داده‌ها می‌توان از داده‌های مربوط به سهام، امور بانکی، اطلاعات وام‌ها، کارت‌های اعتباری، اطلاعات کارت‌های ATM، و کشف کلاهبرداری‌ها می‌باشد.
- 7- ساخت و تولید: حجم زیادی از این داده‌ها روزانه به اشکال مختلفی در کارخانه‌ها تولید می‌شود. ذخیره و دسترسی کارا به این داده‌ها و تحلیل آنها برای صنعت تولید بسیار بااهمیت است.
- 8- کسب و کار و بازاریابی: داده لازم است برای پیش‌بینی فروش، طراحی کسب و کار، رفتار بازاریابی، و غیره.
- 9- شبکه راه‌دور: انواع مختلفی از داده‌ها در این صنعت تولید و ذخیره می‌شوند. آنها برای تحلیل الگوهای مکالمات، دنبال کردن تماس‌ها، مدیریت شبکه، کنترل تراکم، کنترل خطا و غیره، استفاده می‌شوند.
- 10- حوزه علوم: این حوزه شامل مشاهدات نجومی، داده زیستی، داده ژنومیک، و غیره است.

11-WWW: یک حجم وسیع از انواع مختلف داده که در هر جایی از اینترنت پخش شده‌اند.

در بیشتر این حوزه‌ها، تحلیل داده‌ها یک روال دستی بود. یک تحلیلگر کسی بود که با داده‌ها بسیار آشنا بود و با کمک روش‌های آماری، خلاصه‌هایی تهیه و گزارشاتی را تولید می‌کرد. در یک حالت پیشرفته‌تر، از یک پردازنده پیچیده پرسش استفاده می‌شد. اما این روش‌ها با افزایش حجم داده‌ها کاملاً بلااستفاده شدند.

واژه‌های «داده‌کاوی» و «کشف دانش در پایگاه داده» اغلب به صورت مترادف یکدیگر مورد استفاده قرار می‌گیرند. کشف دانش به عنوان یک فرآیند در شکل 1 نشان داده شده است.

کشف دانش در پایگاه داده فرایند شناسایی درست، ساده، مفید، و نهایتاً گوها و مدل‌های قابل فهم در داده‌ها می‌باشد. داده‌کاوی، مرحله‌ای از فرایند کشف دانش می‌باشد و شامل الگوریتم‌های مخصوص داده‌کاوی است، بطوریکه، تحت محدودیتهای مؤثر محاسباتی قابل قبول، الگوها و یا مدل‌ها را در داده کشف می‌کند. بیان ساده‌تر، داده‌کاوی به فرایند استخراج دانش ناشناخته، درست، و بالقوه مفید از داده اطلاق می‌شود. تعریف دیگر اینست که، داده‌کاوی گونه‌ای از تکنیکها برای شناسایی اطلاعات و یا دانش تصمیم‌گیری از قطعات داده می‌باشد، به نحوی که با استخراج آنها، در حوزه‌های تصمیم‌گیری، پیش‌بینی، پیشگویی، و تخمین مورد استفاده قرار گیرند. داده‌ها اغلب حجیم، اما بدون ارزش می‌باشند، داده به تنهایی قابل استفاده نیست، بلکه دانش نهفته در داده‌ها قابل استفاده می‌باشد. به این دلیل اغلب به داده‌کاوی، تحلیل داده‌ای ثانویه گفته می‌شود.

استخراج دانش در پایگاه داده (KDD)، بعنوان روالی برای شناسایی الگوهای معتبر، جدید، بالقوه مفید، و سرانجام قابل فهم در داده‌ها، تعریف شده است. روال سراسری شامل تبدیل داده سطح-پایین به دانش سطح-بالاست. این روال یک روال تعاملی و تکراری است که شامل مراحل زیر می‌باشد. **Error! Reference source not found.**

- 1- درک دامنه کاربرد: این شامل دانش قبلی مرتبط و اهداف کاربرد است.
- 2- استخراج مجموعه داده هدف: این چیزی نیست جز انتخاب یک مجموعه داده یا یک زیرمجموعه از متغیرها، با استفاده از تکنیک‌های رتبه‌بندی و انتخاب است.
- 3- پیش پردازش داده: این مرحله برای افزایش کیفیت داده بکار گرفته شده برای داده‌کاوی، لازم است. همچنین برای بهبود کارایی کاوش داده لازم است. پیش پردازش داده شامل پاکسازی داده، انتقال داده، یکپارچه سازی داده، کاهش یا فشرده‌سازی داده برای نمایش فشرده، و غیره است.
- 4- داده‌کاوی: این مرحله شمل اعمال یکی از الگوریتم‌های داده‌کاوی است.

- 5- تفسیر: شامل تفسیر الگوهای استخراج شده، و تا حد امکان، بصری سازی این الگوهاست. بصری سازی یک کمک کننده مهم در قابل فهم سازی الگوهاست.
- 6- استفاده از دانش استخراج شده: این مرحله شامل تلفیق این دانش با کارایی سیستم و گرفتن تصمیمات عملی براساس این دانش است.
- بیشتر تکنیکهای داده کاوی حداقل به عنوان الگوریتمهای آکادمیک از سالها یا دهه های قبل وجود داشته اند. تنها در دهه اخیر است که داده کاوی تجاری نقش عمده ای را بازی کرده است.

چرا امروزه ما به داده کاوی گرایش داریم؟ در زیر تعدادی از دلایل آن آورده شده: **Error! Reference source not found.**

1. مقدار زیاد داده در دسترس: در دهه اخیر، قیمت سخت افزار بویژه فضای دیسک سخت، بسیار کاهش یافته است. و به دنبال آن، تشکیلات اقتصادی مقدار زیادی از داده ها را از کاربردهای زیادی گردآوری کرد. با این انفجار داده ها، تشکیلات اقتصادی می خواهند که الگوهای پنهان در این داده ها را برای هدایت استراتژی های تجارت خود بکار گیرند. داده کاوی هنگامی بیشترین معنی را پیدا می کند که داده های زیادی وجود داشته باشد. اغلب الگوریتم های داده کاوی نیازمند میزان زیادی از داده ها هستند تا مدلهایی را ترتیب دهند که بعداً برای دسته بندی، تخمین، پیش بینی یا سایر کارکردهای داده کاوی مورد استفاده قرار گیرند.
2. افزایش رقابت: رقابت بعلت وجود بازارهای مدرن و کانالهای توزیع مثل اینترنت و ارتباطات راه دور، بطور فزاینده ای در حال افزایش است. تشکیلات اقتصادی با رقابتهای جهان وب مواجه اند و کلید موفقیت در تجارت، حفظ مشتریان کنونی و بدست آوردن مشتریان جدید است. داده کاوی، تکنولوژی هایی دارد که اجازه می دهد که تشکیلات تجاری فاکتورهایی را برای مواجهه با این زمینه ها تحلیل کند.
3. آماده بودن تکنولوژی آن: داده کاوی قبلاً فقط در حوزه آکادمیک قرار داشت، اما در حال حاضر بسیاری از این تکنولوژی ها کامل شده اند و برای اعمال در صنعت آماده اند. الگوریتم ها، بسیار دقیق تر و کارا تر شده اند و می توانند بطور فزاینده ای داده های پیچیده را مدیریت کنند. بعلاوه رابط برنامه نویسی کاربردهای داده کاوی (APIها)، اکنون استاندارد شده اند، که به توسعه دهندگان این امکان را می دهند که کاربردهای داده کاوی بهتری بسازند.

4. علاقه به مدیریت روابط با مشتریان فراوان است: در طیف وسیعی از صنایع، شرکتها به این بینش رسیده اند که مشتریان برای سازمان حیاتی هستند. و اطلاعات درباره آن مشتریان یکی از دارایی های اساسی سازمان می باشد. اطلاعاتی که شرکتها درباره مشتریانشان دارند نه تنها برای خودشان بلکه برای دیگران هم ارزشمند است. اطلاعات یک محصول است. یک شرکت کارت اعتباری چیزهایی می داند که شرکتهای خطوط هوایی دوست دارند بدانند یعنی چه کسی بلیطهای پرواز متعددی می خرد. گوگل می داند مردم در وب دنبال چه چیزی هستند و از این شناخت با فروش لینکهایی با پشتیبان مالی بهره میبرد. در واقع هر شرکتی که داده های با ارزش جمع آوری می کند در موقعیت یک واسطه اطلاعات قرار دارد.

۱-۲ تاریخچه

با رشد فناوری اطلاعات و روش های تولید و جمع آوری داده ها ، پایگاه داده های مربوط به داده های تبادلات تجاری ، کشاورزی ، اینترنت ، جزییات مکالماتی تلفنی ، داده های پزشکی و ... سریع تر از هر روز جمع آوری و انبارش میشود . لذا از اواخر دهه ی 80 میلادی بشر به فکر دستیابی به اطلاعات نهفته در این پایگاه های داده های حجیم افتاد . زیرا سیستم های سنتی قادر به این کار نبودند . به دلیل رقابت در عرصه های سیاسی ، نظامی ، اقتصادی ، علمی و اهمیت دست یابی به اطلاعات در کم ترین زمان بدون دخالت انسان علم تجزیه و تحلیل داده ها یا داده کاوی پا به عرصه گذاشت .

داده کاوی فر آیندی است که در آغاز دهه 90 مطرح شد و با نگرشی نو به مساله استخراج اطلاعات از پایگاه داده ها میپردازد . از سال 1995 داده کاوی به صورت جدی وارد مباحث آمار شد و در سال 1996 اولین شماره مجله کشف دانش و معرفت از پایگاه داده ها منتشر شد . محقق های نظیر براچمن و آناند کلیه مراحل واقع گرایانه و رو به جلو کشف دانش از پایگاه داده ها را تشخیص دادند .

در حاضر ، داده کاوی مهمترین فناوری جهت بهره برداری موثر از داده های حجیم است و اهمیت آن رو به فزونی است . به طوری که تخمین زده شده است که مقدار داده ها در جهان در هر 20 ماه حدود دو برابر میشود . در یک تحقیق که بر روی گروه های تجاری بسیار بزرگ در جمع آوری داده ها صورت گرفت مشخص گرفت که 19 درصد از این گروه ها دارای پایگاه داده هایی با سطح بیشتر از 50 گیگابایت میباشند و 59 درصد آن ها انتظار دارند که در آینده نزدیک در چنین سطحی قرار گیرند.

در صنایعی مانند کارت های اعتباری و ارتباطات و فروشگاه های زنجیره ای و خرید های الکترونیکی و اسکنر های بار کد خوان هر روزه داده های زیادی تولید و ذخیره میشوند . افزایش سرعت کامپیوتر باعث بوجود آمدن الگوریتم هایی شده است که قدرت تجزیه و تحلیل های بیشتری دارد بدون اینکه محدودیت در زمینه ظرفیت و سرعت کامپیوتر داشته باشد .

در سال 1989 و 1991 کارگاه های کشف دانش و معرفت از پایگاه داده ها توسط پیاتتسکی و همکارانش برگزار شد . در فواصل سال های 1991 تا 1994 کارگاه های کشف دانش و معرفت از پایگاه داده ی توسط فییاد و پیاتتسکی و دیگران برگزار شد . به طور رسمی اصطلاح داده کاوی برای اولین بار توسط فییاد در اولین کنفرانس بین المللی " کشف دانش و داده کاوی " در سال 1995 مطرح شد . امروزه کنفرانس های مختلفی در این زمینه در سراسر دنیا برگزار شد. افزایش داده های بسیار باعث پیدایش فرصت های تازه برای کار در علوم مهندسی و کسب و کار شده است . زمینه داده کاوی و کشف دانش از پایگاه داده ها به عنوان یک رشته علمی جدید در مهندسی علوم کامپیوتر ظهور کرده است . مهندسی صنایع با حوزه های گوناگون و در برداشتن فرصت های بی نظیر اکنون برای کاربرد داده کاوی و کشف دانش از پایگاه داده ها و برای توسعه مفاهیم و روش های تازه در این زمینه آماده ه است. فرآیند های صنعتی زیادی اکنون برای مطمئن شدن از کیفیت سفارشات محصول و کاهش هزینه های محصول به طور خودکار و کامپیوتری شده اند.

فصل دوم

داده کاوی

۱-۲ تعاریف داده کاوی

در متون آکادمیک تعاریف گوناگونی برای داده کاوی ارائه شده اند. در برخی از این تعاریف داده کاوی در حد ابزاری که کاربران را قادر به ارتباط مستقیم با حجم عظیم داده ها می سازد معرفی گردیده است و در برخی دیگر، تعاریف دقیقتر که در آنها به کاوش در داده ها توجه می شود موجود است.

برخی از این تعاریف عبارتند از :

- 1- داده کاوی عبارت است از فرایند استخراج اطلاعات معتبر، از پیش ناشناخته، قابل فهم و قابل اعتماد از پایگاه داده های بزرگ و استفاده از آن در تصمیم گیری در فعالیتهای تجاری مهم .
- 2- فرایند نیم خودکار تجزیه و تحلیل پایگاه داده های بزرگ به منظور یافتن الگوهای مفید اطلاق می شود.
- 3- داده کاوی یعنی فرایند جستجو در یک پایگاه داده ها برای یافتن الگوهایی میان داده ها .
- 4- داده کاوی یعنی تجزیه و تحلیل مجموعه داده های قابل مشاهده برای یافتن روابط مطمئن بین داده ها.
- 5- داده کاوی یعنی استخراج دانش کلان ، قابل استناد و جدید از پایگاه داده های بزرگ.

همانگونه که در تعاریف گوناگون داده کاوی مشاهده می شود، تقریباً در تمامی تعاریف به مفاهیمی چون استخراج دانش ، تحلیل و یافتن الگوی بین داده ها اشاره شده است.

داده کاوی تحلیل اکتشافی داده درمیان چیزهای دیگر نامیده می شود . حجم داده های تولید شده از ثبت پول ، اسکن ، پایگاه داده های خاص موضوعی کل شرکت کشف و تحلیل شده ، کاهش یافته و دوباره استفاده می شوند . تحقیقات درمیان مدل های مختلف پیشنهاد شده برای پیشگویی فروش ، واکنش بازار و سود انجام می شوند . روش های اماری کلاسیک ، اصول داده کاوی می باشند . متد های AI خود کار نیز استفاده می شوند . به علاوه شناسایی اصولی به واسطه متدهای اماری کلاسیک هنوز پایه و اساس داده کاوی هستند . بعضی از ابزارهای توسعه توسط زمینه

تحلیل اماری از طریق کنترل خودکار (همراه با بعضی راهنمایی های کلیدی بشر) در برخورد با داده ها مهار شده اند .

داده کاوی عبارت است از فرایند خودکار کشف دانش و اطلاعات از پایگاه های داده ای . این فرایند تکنیکهایی از هوش مصنوعی را بر روی مقادیر زیادی داده اعمال می کند تا روندها ، الگوها و روابط مخفی را کشف کند . ابزارهای داده کاوی برای کشف دانش یا اطلاعات از داده ها به کاربر اتکا نمی کنند ، بلکه فرایند پیشگویی واقعیت ها را خودکار می سازند . این تکنولوژی نوظهور ، اخیراً به طور فزاینده ای در تحلیل ها مورد استفاده قرار می گیرد .

۱-۱-۲ تعریف جامع:

"داده کاوی فرآیندی است که طی آن با استفاده از ابزار های تحلیل داده به دنبال کشف الگوها و ارتباطات میان داده های موجود که ممکن است منجر به استخراج اطلاعات جدیدی از پایگاه داده گردند، می باشد."

در تعاریفی که از داده کاوی ارائه شد به اصطلاح "فرایند" اشاره شد. حتی در بعضی محیط های حرفه ای این نظر وجود دارد که داده کاوی شامل انتخاب و بکارگیری ابزارهای مبتنی بر کامپیوتر برای حل مسائل فعلی و بدست آوردن یک راه حل بطور اتوماتیک و خودکار میباشد.

برای آموزش داده کاوی، باید بر مفاهیم و روش های اعمال شده برخلاف همه جاذبه های ابزارهای مبتنی بر کامپیوتر که امور را با جزئیات و دستورات با فرمت های خاصی باید به خیلی از سوالات از جمله چگونگی طراحی و استفاده از فرایندها را پاسخ داد به جای بیان جزئیات عملی ابزار مختلف داده کاوی تکیه نمود

2-1-2 مراکز داده چیست

مراکز داده یا سرور فارم ، ساختمان های خاصی هستند که در آن ها انبوهی از تجهیزات آی تی نظیر سیستم های کامپیوتر و مولفه های مرتبط با آن نسب شده تا پیش نیاز های لازم برای عملکرد درست و بی نقض سایت های سازمان فراهم شود . این مرکز همچنین وظیفه ی برقراری نظم و امنیت سایت و تهیه بستر ارتباط سازمان را بر عهده دارد .

2-1-3 برخی از کاربردهای داده کاوی در محیطهای واقعی عبارتند از :

- 1- خرده فروشی
 - 2- تعیین الگوهای خرید مشتریان
 - 3- تجزیه و تحلیل سبد خرید بازار
 - 4- پیشگویی میزان خرید مشتریان از طریق پست (فروش الکترونیکی)
 - 5- بانکداری
 - 6- پیش بینی الگوهای کلاهبرداری از طریق کارتهای اعتباری
 - 7- تشخیص مشتریان ثابت
 - 8- تعیین میزان استفاده از کارتهای اعتباری بر اساس گروههای اجتماعی
 - 9- بیمه
 - 10- تجزیه و تحلیل دعاوی
 - 11- پیشگویی میزان خرید بیمه نامه های جدید توسط مشتریان
 - 12- پزشکی
 - 13- تعیین نوع رفتار با بیماران و پیشگویی میزان موفقیت اعمال جراحی
 - 14- تعیین میزان موفقیت روشهای درمانی در برخورد با بیماریهای سخت
- 2-1-4 مسائل کسب و کار برای داده کاوی. **Error! Reference source not found.**

تکنیکهای داده کاوی می تواند برای کاربردهای بسیاری بکار رود، در زیر تعدادی از مسائل کسب و کار معمولی را که می توان با داده کاوی برای آنها پاسخی یافت، شرح داده می شود:

1. تحلیل رویگردانی: کدام مشتریان بیشتر تمایل دارند بسمت رقیب ما بروند؟ امروزه صنایع تله کام، بانکداری و بیمه، بیش از دیگران در معرض این رقابت ها هستند. بطور متوسط، هر مشترک جدید تلفن همراه، برای شرکت تلفن، هزینه ای بالغ بر 200 دلار در بازار سرمایه گذاری دارد. هر کسب و کاری می خواهد مشتریان بیشتری را کسب کند. تحلیل رویگردانی می تواند به مدیران بازاریابی کمک کند تا دلایل رویگردانی مشتری را درک کند و روابط مشتری را بهبود دهد و وفاداری مشتری را افزایش دهد.
2. فروش مقاطع: مشتریان چه محصولاتی را دوست دارند بخرند؟ فروش مقاطع برای خرده فروشان یک چالش تجاری بزرگ است. بسیاری از خرده فروشان، بویژه خرده فروشان online، برای افزایش فروش خود از این ویژگی استفاده می کنند. برای مثال، اگر شما به یک کتابفروشی Online مثل Amazon.com برای خرید یک کتاب، بروید، شما متوجه شده اید که وب سایت مجموعه ای از پیشنهادات درباره کتابهای مرتبط را به شما پیشنهاد می کند. این پیشنهادات را می توان از تحلیل داده کاوی گرفت.
3. کشف قلب: آیا این ادعای بیمه، کلاهبرداری است؟ شرکت های بیمه، روزانه هزاران دعوی بیمه دارند. برای آنها مهم است که درباره هر مورد تحقیق و بررسی کنند. داده کاوی می تواند برای شناسایی دعاوی ای که بیشتر نادرست هستند، کمک کند.
4. مدیریت ریسک: آیا باید این درخواست وام برای این مشتری تصویب شود؟ این سوال بسیار رایج در سناریوهای بانکی است. تکنیک های داده کاوی می توانند برای رتبه بندی سطح ریسک یک مشتری، بکار روند، و به مدیر در گرفتن یک تصمیم مناسب برای هر کاربرد، کمک کنند.
5. قطعه بندی مشتریان: چه کسی مشتری من است؟ قطعه بندی مشتریان به مدیران بازاریابی کمک می کند که تفاوت های پروفایل های مشتریان را درک کنند و عمل بازاریابی مناسبی را بر مبنای هر بخش، انجام دهند.
6. تبلیغات هدف دار: چه بنر تبلیغی باید برای یک بازدیدکننده خاص، نمایش داده شود؟ فروشندگان وب و سایت های پورتال تمایل دارند که محتوای خود را برای مشتریان سفارشی کنند. با استفاده از الگوهای ناوبری مشتری یا خرید online، این سایت ها می توانند راه حل های داده کاوی را برای نمایش تبلیغات هدف دار برای مشتریان، بکار برند.
7. پیش بینی فروش: من چه نمونه هایی را در این فروشگاه در این هفته خواهم فروخت؟ تکنیک های پیش بینی داده کاوی می تواند برای پاسخ به انواع این پرسش های مرتبط با زمان، بکار روند.

۲-۲ چرخه تعالی داده کاوی چیست؟

باید بتوان داده‌ها را به اطلاعات، اطلاعات را به عمل و عمل را به ارزش تبدیل کرد. این را در یک کلام چرخه تعالی داده‌کاوی می‌نامند. به منظور دستیابی به این هدف لازم است داده‌کاوی به صورت یک فرآیند ضروری در کنار سایر فرآیندها نظیر بازاریابی، فروش، پشتیبانی مشتری و کنترل موجودی درآید.

در ادبیات بازاریابی داده‌کاوی را بسیار آسان جلوه می‌دهند بطوریکه تنها کفایت الگوریتم‌های خودکار تهیه شده توسط بهترین دانشمندان دانشگاهی نظیر شبکه‌های عصبی، درخت‌های تصمیم‌گیری و الگوریتم ژنتیک را به کار برد تا در مسیر موفقیت قرار گیرید. اگرچه این الگوریتم‌ها مهمند، اما راه حل‌های داده‌کاوی چیزی فراتر از مجموعه‌ای از تکنیکها و ساختارهای داده‌ای قوی است. این تکنیکها را باید در جای مناسب و با داده‌های صحیح بکار برد. چرخه تعالی داده‌کاوی یک فرآیند یادگیری تکراری و مرحله‌ای است که بر اساس نتایج در طول زمان تهیه می‌گردد. موفقیت در بکارگیری داده‌ها، وضعیت یک سازمان را از واکنشی به کنشگرا تبدیل خواهد نمود. با استفاده از چرخه تعالی داده‌کاوی مطرح شده در این کارگاه، بیشترین سود از تکنیکهای داده‌کاوی به دست می‌آید.

این چرخه شامل چهار مرحله است:

1. تشخیص مشکل کسب و کار و تجارت
 2. واکاوی داده‌ها برای تبدیل آن‌ها به اطلاعات عملی
 3. کار روی اطلاعات
 4. بررسی نتایج
- کلید موفقیت، در گنجاندن داده کاوی در فرآیندهای تجاری است و اینکه بتوان راههای ارتباطی میان داده کاوان تکنیکی و کاربران تجاری نتایج یافت.

1- تعیین فرصت‌های کسب و کار و تجارت :

چرخه تعالی داده کاوی با تعیین فرصت‌های واقعی کسب و کار و تجارت آغاز می‌گردد. متأسفانه تحلیل‌گران زیادی هستند که معمولاً تلاش‌هایشان بی‌ثمر می‌ماند زیرا آنها مسائلی را حل می‌کنند که به

کسب و کار و تجارت کمکی نمی نماید. داده کاوان خوب می خواهند این وضعیت برایشان پیش نیاید. اجتناب از تلاشهای تحلیلی بی حاصل با اراده برای کار روی نتایج آغاز می گردد.

بیشتر فرایندهای تجاری معمول کاندیدهای خوبی برای داده کاوی هستند:

1- برنامه ریزی برای معرفی محصول جدید

2- برنامه ریزی برای بازاریابی مستقیم

3- فهم احساسات و گلایه های مشتریان

اگر نتوان نتایج داده کاوی را سنجید نمی توان از آن چیزی یاد گرفت و آنگاه هیچ چرخه تعالی وجود نخواهد داشت. سنجش تلاش های صورت گرفته در گذشته سوالاتی را درباره تجارت، فرصت های جدیدی از داده کاوی را فراهم می کند:

1- چه نوع مشتریانی به جدید ترین اقدام پاسخ داده اند؟

2- بهترین مشتریان کجا زندگی می کنند؟

3- آیا صرف زمان طولانی در پای دستگاه های خودپرداز دلیل گلایه آنهاست؟

4- آیا مشتریان سودآور از خدمات پشتیبانی مشتریان استفاده می کنند؟

5- چه نوع محصولات دیگری را باید با یک محصول خاص تبلیغ نمود؟

مصاحبه با خبرگان تجارت، روش خوب دیگری برای آغاز است؛ زیرا ممکن است کسانی که در داخل تجارت هستند، با داده کاوی آشنا نباشند. آنها ممکن است ندانند چگونه روی نتایج کار کنند.

با بیان ارزش داده کاوی برای یک سازمان، چنین مصاحبه هایی امکان ایجاد ارتباطی دوسویه را فراهم می کند

2- مشکلات در راه کسب توانایی استفاده از نتایج داده کاوی :

داده ها در اشکال گوناگون و از سیستم های چندگانه ای ظاهر می شوند. تشخیص منابع درست داده ها و گرد آوردن آنها کنار یکدیگر از عوامل حیاتی در موفقیت می باشد.

هر پروژه داده کاوی مسائل خاص خود در ارتباط با داده ها را دارد: سیستم های اطلاعاتی ناهماهنگ، اطلاعاتی که هر چند ماه یکبار باز نویسی می شوند. در اینجا است که الگوریتم های داده کاوی بکار می آیند.

3- انجام عمل :

انجام عمل، هدف چرخه متعالی داده کاوی است. همانطور که ذکر شد، عمل می تواند اشکال مختلفی داشته باشد. داده کاوی سبب می شود تصمیمات کسب و کار آگاهانه تر اتخاذ شود. و در طول زمان، تصمیمات با آگاهی بیشتر، نتایج بهتری بدنبال دارد.

اعمال معمولاً به سمت آنچه کسب و کار انجام می دهد، پیش می روند:

1- ارسال پیام ها به مشتریان و انتظارات از طریق پست مستقیم، پست الکترونیک، بازاریابی از راه دور و غیره.

2- با داده کاوی، این امکان بوجود می آید که پیام های مختلف به افراد مختلفی فرستاده شود.

3- اولویت بندی سرویس های مشتریان

4- تعدیل کردن سطوح موجودی

5- و غیره

لازم است که نتایج داده کاوی روال های کسب و کار را تغذیه کند تا مشتریان لمس شوند و روی روابط با مشتریان تاثیر گذاشته شود.

4- اندازه گیری نتایج :

اهمیت اندازه گیری نتایج، در حال حاضر پررنگ خواهد شد. علی رغم اهمیت آن، این مرحله ای است در چرخه تعالی که اغلب چشم پوشی می شود. حتی با اینکه بر ارزش اندازه گیری و بهبود پیوسته، بسیار تاکید شده است، معمولاً به نسبت استحقاق آن، توجه کمتری به آن می شود. چقدر موارد کسب و کاری وجود دارد که پیاده سازی شده اند و هیچ کسی برنگشته تا ببیند آیا آنچه پیاده سازی شده واقعا با طرح سازگار است یا خیر؟ افراد تلاش هایشان را با مقایسه و یادگیری کنترل می کنند، با طرح سوالاتی درباره اینکه چرا طرح ها با آنچه واقعا اتفاق افتاده سازگار است یا سازگار نیست. آنچه برای افراد صدق می کند برای سازمان ها نیز کاربرد دارد.

زمان شروع تفکر درباره اندازه گیری در شروع زمان شناسایی مسئله کسب و کار است. چطور می توان نتایج را اندازه گیری کرد؟ یک معیار مناسب دیگر، سنجش افزایش فروش در مغازه ها و یا مناطق مشخص است. این افزایشها را می توان به تلاشهای خاص بازاریابی مرتبط نمود. ممکن است چنین سنجش هایی از آنجا که نیازمند اطلاعات جزئی فروش است مشکل باشد. با این وجود اگر هدف افزایش فروش است روشی برای بررسی مستقیم آن لازم است.

ایده خوبی است که به هر تلاش داده کاوی به عنوان یک مورد کوچک تجاری نگاه کنیم.

مقایسه انتظارات با نتایج عملی این امکان را ایجاد میکند که فرصت‌های آتی را تشخیص دهیم تا از آنها در دوره‌های بعدی چرخه تعالی بهره ببریم. هر اقدام داده‌کاوی چه موفقیت‌آمیز باشد یا نه، حاوی نکاتی است که می‌تواند برای تلاش‌های آتی بکار رود.

سؤال این است که چه چیزی را بررسی نماییم و چگونه به آن پردازیم، تا بهترین ورودی‌ها را برای استفاده‌های آتی داشته باشیم.

۲-۳ فرایند داده کاوی:

داده کاوی نیازمند شناسایی یک مسئله همراه با مجموعه‌ای از داده‌هاست که می‌تواند منجر به درک بهتر و مدل‌های کامپیوتری به منظور فراهم آوردن تحلیل‌های آماری یا مقاصد دیگر گردد. این کار با ابزارهای تجسم فکری که داده را نمایش می‌دهد یا تحلیل آماری کاملاً بنیادی مثل تحلیل همبستگی ممکن است پشتیبانی شود.

ابزارهای داده کاوی نیاز دارند که متنوع، درجه بندی شده، قادر به واکنش‌های پیشگویی دقیق بین عملیات و نتایج و اجرای خودکار باشند. تنوع به توانایی ابزاری اشاره دارد که انواع وسیعی از مدل‌ها را به کار گیرد. ابزارهای مقیاس اشاره به این موضوع می‌کنند که اگر این ابزارها روی یک مجموعه داده کوچک کار کنند باید روی مجموعه داده‌های بزرگتر نیز کار کنند. کنترل خودکار سودمند است اما کاربردش وابسته به اموری است.

بعضی کارکردهای تحلیلی اغلب خودکار هستند اما تنظیمات بشری مقدم بر اجرای رویه‌ها مورد نیاز هستند. در حقیقت قضاوت تحلیل‌گر انتقادی به اجرای موفقیت‌آمیز داده کاوی است. انتخاب مناسب داده‌ها که شامل تحقیقات می‌باشد حیاتی است. تغییر شکل داده نیز اغلب مورد نیاز است.

متغیرهای بسیاری بازده زیادی تولید می‌کنند. در حالی که تعداد کمی می‌توانند بر روابط کلیدی در داده‌ها مسلط شوند. درک بنیادی مفاهیم آماری برای داده کاوی موفق ضروری است.

داده کاوی به سرعت با فواید زیادی برای تجارت بسط داده می شود . دو ناحیه کاربرد با بیشترین سود ، استفاده از تقسیم مشتری از طریق سازمان های بازاریابی برای تشخیص احتمالات بیشتر حاشیه ای واکنش به شکل های متفاوت رسانه بازاریابی ، و بانک ها با استفاده از داده کاوی برای پیش گویی دقیق تر احتمال این که افراد به پیشنهادهای که سرویس های مختلف ارائه می کنند واکنش می دهند ، می باشد .

بسیاری از شرکت ها از این تکنولوژی برای شناسایی مشتریان استفاده می کنند تا اینکه بتوانند خودشان سرویسی را که نیاز است فراهم آورند تا آن مشتریان را از دست ندهند .

متدلوژی داده کاوی و بهترین تمرین های آن:

در بخش قبل ما چرخه تعالی داده کاوی و مراحل آنرا بررسی کردیم. حالا زمان آن رسیده که به داده کاوی بعنوان یک روال تکنیکی نگاهی بیندازیم. رئوس مطالب در سطح بالا، همان باقی می ماند ولی تاکیدها منتقل می شود. به جای شناسایی مسئله کسب و کار، حالا می خواهیم توجه مان را از مسئله کسب و کار به مسئله داده کاوی منتقل کنیم.

بهترین راه برای دوری از شکست چرخه داده کاوی، درک راه های عدم موفقیت آن و اخذ تدابیر پیشگیرانه است. در طول سالها، مولفان با راه های زیادی برای پروژه های داده کاوی مواجه شده اند که به اشتباه رفته اند. در پاسخ، ما مجموعه ای از عادات کارآمد را توسعه داده ایم - چیزهایی که ما باید برای هموار کردن مسیرمان از جمله ابتدایی مسئله کسب و کار تا یک مدل پایدار که نتایج قابل اندازه گیری و قابل اعمالی را تولید می کند، انجام دهیم. داده کاوی یک روال تکراری طبیعی است. لازم است برخی مراحل را چندین بار تکرار کرد، ولی هیچیک را نمی توان کاملاً نادیده گرفت.

با افزایش پیچیدگی راه حل داده کاوی، نیاز به وجود یک مشی سفت و سخت بیشتر است. پس از اینکه نیاز به وجود متدلوژی با راه های مختلفی توضیح داده شد، با استفاده از پرسش های خاص منظوره نشان داده خواهد شد که در صورت نبود یکی از مراحل آن، تلاش های داده کاوی با شکست مواجه خواهند شد. در پایان، 4 مرحله از چرخه تعالی داده کاوی به 11 مرحله متدلوژی داده کاوی تبدیل می شود.

چرا یک متدلوزی داشته باشیم؟

داده‌کاوی راهی است برای یادگیری از گذشته تا بتوان تصمیمات بهتری در آینده گرفت. داشتن یک متدلوزی برای اجتناب از دو نتیجه نامطلوب روال یادگیری است:

1- یادگیری چیزهایی که درست نیستند

2- یادگیری چیزهایی که درستند ولی کارآمد نیستند

یادگیری چیزهایی که درست نیستند بسیار خطرناک‌تر از یادگیری چیزهایی است که کارآمد نیستند، چون بسیار مهم است که تصمیمات سازمان براساس اطلاعات درست گرفته شود. یافته‌های داده‌کاوی به نظر می‌رسد که قابل اعتماد باشند چون بر اساس داده‌های واقعی در یک حالت ظاهراً علمی هستند. این نمود قابل اعتماد می‌تواند فریبنده باشد. ممکن است داده نادرست باشد یا مرتبط با سوال نباشد. انتقال داده‌ها مانند خلاصه‌سازی‌ها، ممکن است مشکل‌دار باشند یا اطلاعات مهمی را پنهان کنند. در زیر درباره برخی از مسائل بسیار رایجی که منجر به استنتاجات اشتباه می‌شوند، بحث می‌کنیم.

الگوهایی که ممکن است هیچ قانون اصولی را ارائه نکنند :

ارقام اغلب دروغ نمی‌گویند، ولی دروغ‌ها می‌توانند رقم شوند. وقتی می‌خواهیم الگوها را در داده‌ها بیابیم، ارقام واقعا نمی‌خواهند با هدف پیشنهاد چیزهای نادرست، بما دروغ بگویند. راه‌های زیادی وجود دارد که الگوهایی ساخته شود که هر مجموعه داده تصادفی، اگر به اندازه کافی آزمایش شود، یکی را آشکار خواهد کرد. ما انسانها، به الگوها در زندگی مان بسیار وابسته ایم و دوست داریم آنها را همه جا ببینیم حتی اگر وجود نداشته باشند. وقتی به آسمان شب نگاه می‌کنیم و یک مجموعه تصادفی ستارگان را می‌بینیم، آنها را به شکل دب اکبر، دب اصغر و غیره می‌بینیم. گاهی حتی الگوهای مربوط به نجوم را نشانه‌هایی می‌بینیم که از آینده خیر می‌دهند.

احتمالا دلیل اینکه بشر این قبیل وابستگی‌ها را برای الگوها استنتاج می‌کند این است که این الگوها اغلب برخی حقایق را درباره کار جهان، بازتاب می‌دهند.

چالش پیش روی کاونده داده این است که کشف کند که کدام الگوها پیشگویانه هستند و کدام نیستند. به الگوهای زیر دقت کنید، همه آنها در مقالات بعنوان اینکه ارزش پیش بینی دارند، ذکر شده‌اند:

1- حزبی که در ریاست جمهوری را برنده نشده، در انتخابات off-year کرسی‌های کنگره را برنده می‌شود.

2- وقتی که لیگ امریکایی سری‌های جهانی را برنده می‌شود، جمهوری خواهان کاخ سفید را می‌برند.

3- در مناظرات ریاست جمهوری آمریکا، کسی که ق‌بلندتر است، معمولا برنده است.

الگوی اول، به نظر می‌رسد که با عبارات صرف سیاسی قابل توضیح باشد. چون در اینجا یک توضیح متضمن آن وجود دارد، این الگو به نظر می‌رسد به آینده نظر دارد بنابراین ارزش پیشگویانه خواهد داشت. دو عبارت بعدی، یکی که حاوی حوادث ورزشی است، به نظر می‌رسد که به وضوح ارزش پیشگویی نداشته باشد. هیچ اهمیتی وجود ندارد که چند بار در گذشته فاتحان جمهوری خواهان و لیگ آمریکایی، مشترک بوده اند و این چیزی است که تحقیقی برای آن انجام نشده است. در مورد ق‌د کاندیدها چگونه؟ از سال 1945 تا کنون، همیشه کاندیدهای ق‌د بلندتر، برنده انتخابات بوده‌اند. اما به نظر می‌رسد که ق‌د هیچ ربطی به ریاست جمهوری نداشته باشد. از سوی دیگر، ق‌د، ارتباط مثبتی روی درآمد و سایر موفقیت‌های اجتماعی دارد، و انتخاب‌کنندگان، بطور عمدی یا غیرعمدی، کاندیدهای بلندتر را ترجیح می‌دهند.

چیدمان مدل ممکن است بازتاب دهنده جمعیت وابسته نباشد :

چینش مدل، مجموعه ای از داده‌های تاریخی است که برای توسعه مدل داده‌کاوی بکار می‌رود. برای اینکه تفاسیر برگرفته از چینش مدل، معتبر باشد، چیدمان مدل باید منعکس کننده جمعیتی باشد که مدل قصد تشریح، طبقه بندی یا رتبه‌بندی آنرا داشته است. یک مثال که جمعیت والد به خوبی منعکس نشده است، جانبداری است. استفاده از یک نمونه جانبدارانه، بعنوان یک چیدمان

مدل، یک دستورالعمل برای موارد آموزشی ای است که درست نیستند. اجتناب از آن نیز سخت است. فرض کنید:

- 1- مشتریانی که مطابق با انتظارات نیستند.
 - 2- ممیزی پاسخ دهندگان مانند آنها که پاسخ نداده اند نیست.
 - 3- افرادی که Email می خوانند مانند آنها که نمی خوانند نیستند.
 - 4- افرادی که در یک سایت وب ثبت نام می کنند، مانند آنها که موفق به ثبت نام نمی شوند، نیستند.
 - 5- رکوردهایی که هیچ مقدار گم شده ای ندارند، جمعیت متفاوتی را به نسبت رکوردهایی که مقادیر گم شده دارند، منعکس می کنند.
- در انتخاب و نمونه گیری از داده ها برای مدل، بسیار دقیق باشید؛ زیرا در موفقیت داده کاوی بسیار موثر خواهد بود.

یادگیری چیزهایی که درست ولی بلااستفاده اند :

یادگیری چیزهایی که بلااستفاده اند، اگرچه به خطرناکی یادگیری چیزهایی که نادرستند نیست، ولی بسیار رایج است.

1- یادگیری چیزهایی که در حال حاضر می دانیم

داده کاوی باید اطلاعات جدیدی تولید کند. بسیاری از الگوهای قوی در داده، چیزهایی را که می دانیم، ارائه می کنند. مثلا افراد بالای سن بازنشستگی، تمایلی به پاسخ دادن به طرح های ذخیره بازنشستگی، ندارند. افرادی که در جایی زندگی می کنند که برج مخابرات ندارند، تمایلی به خرید تلفن ندارند.

قویترین الگوها، اغلب قانون های کسب و کار را بازتاب می دهند. اگر داده کاوی کشف کند که افرادی که امکان قطع تلفن های نامشخص را دارند، caller ID نیز دارند، این محتمل است چون که این امکان با بسته ای فروخته می شود که شامل caller ID نیز هست. ما بسیاری از این الگوهای کشف شده را در داده کاوی می بینیم. این الگوها نه تنها مطلوب نیستند، که ممکن است قدرت آنها، الگوهای بدیهی دیگری را محو کند.

یادگیری چیزهایی که در حال حاضر می دانیم، یک هدف مفید دارد اینکه ثابت می کند که در سطح تکنیکی، تلاش داده کاوی کار می کند و داده صحیح است. این می تواند کاملاً تسلی بخش باشد. اگر داده و داده کاوی اعمال شده بر آن به اندازه کافی قدرتمند باشد که چیزهای شناخته شده را بدرستی کشف می کند، پس این اطمینان را می دهد که سایر کشفیاتش نیز درست باشد.

2- یادگیری چیزهایی که نمی توان آنها را مورد استفاده قرار داد

این زمانی اتفاق می افتد که داده کاوی روابط پوشش داده نشده ای را که هم درستند و هم قبلاً شناخته نشده اند را کشف می کند ولی مشکل است که از آنها استفاده کرد. مثلاً سابقه اعتباری یک مشتری ممکن است یک دعوی بیمه را در آینده پیش بینی کند، اما تنظیم کننده، اتخاذ تصمیم را بر مبنای آن منع کند.

داده کاوی ممکن است خروجی های دیگری را که خارج از کنترل شرکت هستند را پیش بینی کند. ممکن است یک محصول برای یک آب و هوا مناسب تر از آب و هوای دیگری باشد، ولی سخت است که آب و هوا را کنترل کرد.

نکته: بعضی اوقات ممکن است یک تصور اشتباه باعث شود که اطلاعات جدید را بی استفاده بدانیم. پیش بینی های رویگردانی مشتریان ممکن است دیگر برای بکارگیری حفظ مشتریان کنونی بسیار دیر باشد، اما می تواند ما را برای یافتن راههایی برای تغییر کانال های ارتباطی مان با مشتریان آینده، ترغیب کند.

2-3-1 مدل ها، پروفایل سازی، و پیش بینی :

تکنیک های داده کاوی شرح داده شده در اینجا، همگی برای یادگیری چیزهای جدید با ساخت مدل ها بر اساس داده ها، طراحی شده اند.

سراسر داده کاوی درباره ساختن مدل هاست. همانطور که در شکل 4 نشان داده شده است، مدل ها مجموعه ای ورودی می گیرند و یک خروجی تولید می کنند. داده ای که برای ساختن مدل بکار می رود مجموعه مدل نامیده می شود. و هنگامی که مدل ها روی داده های جدید اعمال می شوند، مجموعه رتبه نامیده می شوند. مجموعه مدل سه جزء دارد:

- 1- مجموعه آموزشی که برای ساختن یک مجموعه مدل بکار می‌رود.
 - 2- مجموعه اعتبارسنجی که برای انتخاب بهترین مدل از میان اینها بکار می‌رود.
 - 3- مجموعه آزمون که برای تعیین چگونگی عملکرد مدل بر روی داده‌های دیده نشده، بکار می‌رود.
- تکنیک‌های داده‌کاوی را می‌توان برای ساختن سه نوع مدل برای سه نوع وظیفه، استفاده کرد: پروفایل‌سازی توصیفی، پروفایل‌سازی جهتدار، و پیش‌بینی.
- مدل‌های توصیفی، تشریح می‌کنند که چه چیزی در داده است. خروجی آنها، نمودارها یا اعداد یا گرافیک‌هاست. از سوی دیگر، پروفایل‌سازی جهتدار و پیش‌بینی، از ساخت مدل، هدفی را دنبال می‌کنند. تفاوت این دو مدل در چارچوب‌های زمانی است که در شکل شان داده شده است. در مدل پروفایل‌سازی، هدف در همان چارچوب زمانی ورودی قرار دارد در حالیکه در مدل پیش‌گویی، هدف در چارچوب زمانی بعدی است.

2-3-2 پروفایل‌سازی :

پروفایل‌سازی یک دستاورد آشنا برای بسیاری از مسائل است. نیازی به هیچ تحلیل داده سطح بالا و پیچیده ای ندارد. پروفایل‌ها، اغلب بر اساس متغیرهای سرشماری، مثل جنسیت، سن و مکان جغرافیایی هستند. از زمانی که تبلیغات بر طبق این متغیرها فروخته می‌شوند، پروفایل‌های سرشماری می‌توانند مسقیما در استراتژی‌های رسانه قرار گیرند. پروفایل‌های ساده ای برای حق بیمه شرکتهای بیمه بکار می‌روند. یک مرد 17 ساله، بنسبت یک زن 60 ساله، حق بیمه بیشتری برای ماشینش باید پرداخت کند. همچنین کاربردهایی که برای پرداخت حق بیمه عمر هستند، سوالاتی درباره سن، جنسیت و سیگاری بودن می‌پرسند.

علی‌رغم این قدرت، پروفایلینگ محدودیت‌های جدی‌ای دارد. اول اینکه قادر به تشخیص انگیزه و نتیجه نیست. مادامی که پروفایل‌سازی مبتنی بر متغیرهای سرشماری آشنا باشد، این نمی‌تواند قابل توجه باشد.

با داده‌های رفتاری، مسیر علیت همیشه روشن نیست. یک جفت مثال واقعی از پروژه‌های داده‌کاوی واقعی را فرض کنید:

1- افرادی که با سپرده‌های اعتباری شان خرید کرده اند، در حساب‌های ذخیره‌شان، مقدار کم یا هیچ پولی ندارند.

2- مشتریانی که از پست صوتی استفاده می‌کنند، تماس‌های کوتاه زیادی با شماره‌شان می‌گیرند.

نگه داشتن پول در حساب‌های ذخیره، رفتار رایج دارندگان سپرده‌های اعتباری است. مثل مرد بودن که خصوصیت رایج نوشندگان نوشابه است. شرکت‌های نوشابه، مشتریان مرد را برای بازار محصولشان جستجو می‌کنند، همینطور بانک‌ها، برای فروش سپرده‌های اعتباری‌شان، کسانی را حساب ذخیره خالی دارند، جستجو می‌کنند؟ شاید نه! احتمالاً دارندگان سپرده‌های اعتباری هیچ پولی در حساب‌های ذخیره خود ندارند چون پولشان را برای خرید از اعتبارشان استفاده می‌کنند. یک دلیل رایج برای نداشتن پول در حساب‌های ذخیره، نداشتن هیچ پولی است و کسانی که پولی ندارند برای خریدن کارت اعتباری اصلاً مطلوب نیستند.

2-1-8 پیش بینی :

پروفایل سازی، داده‌های گذشته را برای تشریح آنچه که در گذشته رخ داده است بکار می‌برد. پیش‌بینی یک گام به جلوست. و داده‌های گذشته را برای پیش‌بینی آنچه که در آینده رخ خواهد داد، بکار می‌برد. این استفاده از داده، قدرت بیشتری دارد.

ساختن مدل پیشگویانه، نیاز به جدایی زمانی بین ورودی‌های مدل یا پیشگو، و خروجی مدل، چیزی که پیش‌بینی شده، دارد. اگر این جدایی حاصل نشود، مدل کار نخواهد کرد. این مثالی است از اینکه چرا مهم است که از متدلوژی داده‌کاوی پیروی کنیم.

2-4 متدلوژی :

متدلوژی داده‌کاوی 11 مرحله دارد:

1- تبدیل مسئله کسب و کار به مسئله داده‌کاوی

- 2- انتخاب داده مناسب
- 3- رسیدن به شناخت داده
- 4- ایجاد یک مجموعه مدل
- 5- تثبیت مسئله با داده
- 6- تبدیل داده برای آوردن اطلاعات به سطح
- 7- ساختن مدل ها
- 8- ارزیابی مدل ها
- 9- استقرار مدلها
- 10- ارزیابی نتایج
- 11- شروع دوباره

مرحله 1: تبدیل مسئله کسب و کار به مسئله داده‌کاوی :

اهداف داده‌کاوی برای یک پروژه خاص نباید با عبارت کلی نوشته شده باشد، مثلاً:

- 1- کسب درکی از رفتار مشتری
- 2- کشف داده‌های معنی‌دار از داده‌ها
- 3- یادگیری چیزهای جالب

تمام اینها اهداف شایسته‌ای هستند ولی حتی هنگامی که آنها بدست آیند، اندازه‌گیری آنها مشکل است. پروژه‌هایی که اندازه‌گیری آنها سخت است، ارزش‌گذاری آنها سخت است. تا آنجا که ممکن است، باید اهداف کلی را به اهداف خاص دیگری شکسته شوند که نظارت بر روال دسترسی به آنها آسانتر شود. بدست آوردن درکی از رفتار مشتری، می‌تواند به اهداف زیر شکسته شود:

- 1- شناسایی مشتریانی که دوست ندارند اشتراکشان را تجدید کنند.
 - 2- طراحی یک طرح مکالمه‌ای که رویگردانی مشتریان خانگی را کاهش دهد
 - 3- رتبه‌بندی تمام مشتریانی که گرایش به اسکی دارند.
- این اهداف واقعی، فقط برای نظارت آسانتر نیستند بلکه تبدیل آنها به مسائل داده‌کاوی نیز آسانتر است.

برای تبدیل یک مسئله کسب و کار به یک مسئله داده‌کاوی، باید آنرا بصورت یکی از وظایف داده‌کاوی شرح داده شده دربخش قبل، دوباره فرمول سازی کرد

از نتایج چگونه استفاده می‌شود؟

این یکی از سوالات بسیار مهمی است که در این مرحله باید از خود پرسیم. بطور شگفت‌انگیزی، پاسخ در ابتدا این است: مطمئن نیستم! اما پاسخ به این سوال بسیار مهم است چون کاربردهای مختلف، راه‌حل‌های مختلفی را نیز پیشنهاد می‌کنند.

نتایج چگونه تحویل داده می‌شوند؟

پروژه‌های داده‌کاوی مختلف، ممکن است چندین نوع نتایج قابل تحویل داشته باشند. مثلاً وقتی هدف اوله پروژه، بدست آوردن درکی از مشتری باشد، نتیجه قابل تحویل اغلب یک گزارش، چارت یا نمودار یا گراف است. شکل قابل تحویل می‌تواند نتایج داده‌کاوی را تحت تاثیر قرار دهد.

نقش کاربران کسب و کار و تکنولوژی اطلاعات :

تنها راه پاسخ دادن مطلوب به سوالات مطرح شده در بالا، درگیر کردن صاحبان کسب و کار در کشف چگونگی استفاده از نتایج داده‌کاوی، و درگیر کردن کارکنان فن‌آوری اطلاعات در کشف نحوه تحویل این نتایج است.

مرحله 2: انتخاب داده مناسب :

داده‌کاوی به داده نیاز دارد. در بهترین حالت ممکن، داده مورد نیاز در یک انبار داده مجتمع، پالایش شده، در دسترس، با سابقه درست، و بطور متناوب در حال اصلاح، قرار دارد. اما در واقعیت تمام این موارد ممکن نیست. منابع داده مفید و دردسترس، از مسئله‌ای به مسئله دیگر، و از صنعتی به صنعت دیگر، متنوع هستند.

چه چیزی در دسترس است ؟

اولین جایی که باید بدنبال داده گشت، یک انبار داده مجتمع است. داده‌ها در انبار، پالایش، اعتبارسنجی و از چندین منبع باهم گردآوری شده‌اند. یک مدل داده تکی، این اطمینان را می‌دهد که فیلدها بطور مشابهی نامگذاری شده‌اند، معنای یکسانی، و انواع داده سازگاری را در پایگاه داده دارند. یک انبار داده یکی شده، یک مخزن سابقه ای است؛ داده‌های جدید به آن اضافه می‌شوند ولی داده‌های قبلی تغییر داده نمی‌شوند. از آنجایی که برای پشتیبانی تصمیم‌گیری طراحی شده‌اند، انبار داده، داده‌های با جزئیاتی را که برای داده‌کاوی در سطح درستی پذیرفته شود، تهیه می‌کند. تنها مسئله این است که در بسیاری از سازمان‌ها، واقعا چنین انبار داده ای وجود ندارد یا اینکه یک یا چند انبار داده وجود دارد ولی در سطح قابل قبولی نیستند. در شروع این مورد، داده‌کاو باید داده‌ها را از میان پایگاه داده‌های مختلف اداری و از دل سیستم‌های عملیاتی مختلف بیرون بکشد. ممکن است سیاست‌های معنی‌دار و تلاش‌های برنامه نویسی‌ای برای گرفتن داده بشکل کارآمد برای کشف دانش از چنین سیستم‌هایی نیاز شود. در برخی موارد، روال‌های عملیاتی برای تهیه داده، تغییر خواهند کرد.

چه مقدار داده کافی است ؟

متأسفانه پاسخ ساده ای برای این سوال وجود ندارد. پاسخ به الگوریتم‌های خاصی که بکار گرفته می‌شود و پیچیدگی داده، بستگی دارد. در مواردی که داده کم است، داده‌کاوی نه تنها کم اثر است بلکه بلااستفاده نیز هست. داده‌کاوی بیشترین فایده را زمانی دارد که حجم خالصی از داده، الگوهای را پنهان کند که قابل جستجو در پایگاه داده‌های کوچک هستند.

در داده‌کاوی، داده‌ها هر چه بیشتر، بهتر است ولی باید به دو نکته توجه کرد. اولین نکته، رابطه بین اندازه مجموعه مدل و تراکم آن است. منظور از تراکم پخش اثرات مطلوب‌هاست. اغلب متغیر هدف، چیزی را که نسبتا نادر است، ارائه می‌کند. برای دارندگان کارت‌های اعتباری، ارتکاب تقلب نادر است. اینکه مشتریان یک روزنامه، اشتراکشان را لغو کنند، کم پیش می‌آید. برای مجموعه مدل، مطلوب این است که با اعداد برابر با هر یک از خروجی‌های روال ساخت

مدل، متوازن باشد. یک مجموعه کوچکتر متوازن به مجموعه بزرگتری که موارد نادر را داشته باشد، ترجیح داده می‌شود.

دومین نکته، این است که وقتی مجموعه مدل برای ساختن مدل های خوب و پایدار، به اندازه کافی بزرگ است، بزرگتر ساختن آن خنثی کننده تولید است، چون هر چیزی برای اجرا روی مجموعه داده، زمان بیشتری می‌برد. از آنجا که روال داده‌کاوی یک فرایند تکراری است، اگر هر اجرا از یک روتین مدل‌سازی به جای چند دقیقه، ساعت‌ها وقت بگیرد، زمان صرف شده برای رسیدن به نتایج ممکن است طولانی شود.

یک تست ساده برای اینکه بدانیم مجموعه داده‌ها برای شروع به اندازه کافی مناسب است این است که اندازه نمونه را دوبرابر کنیم و بهبود درستی مدل را اندازه‌گیری کنیم. اگر مدل ساخته شده از داده‌های بزرگتر، بطرز معنی‌داری بهتر از مدل کوچکتر بود، بدان معنی است که مجموعه نمونه ما به اندازه کافی بزرگ نیست. اما اگر بهبود چشمگیری در مدل حاصل نشد، احتمالاً مدل اصلی مناسب خواهد بود.

چقدر سابقه لازم است؟

داده‌کاوی، داده‌های گذشته را برای پیش‌بینی آینده بکار می‌برد. ولی چه مقدار از سابقه را باید با داده‌ها آورد؟ این یک سوال ساده دیگری است که پاسخ ساده‌ای ندارد. اولین چیز برای توجه فصلی بودن است. برخی کسب و کارها درجه‌ای از فصلی بودن را دارند. مثلاً سفرهای فراغت در تابستان اتفاق می‌افتد. فروش‌ها بیشتر در یک چهارم پایانی سال انجام می‌شود. باید به اندازه کافی داده قبلی برای رویدادهای دوره‌ای از این دست، گرفته شود.

از سوی دیگر، داده‌های گذشته‌های دور بدلیل تغییر شرایط بازار، نمی‌توانند خیلی مفید باشند. برای بیشتر کاربردهای متمرکز بر مشتری، سوابق دو تا سه سال، مناسب است.

چند تا متغیر؟

داده‌کاوه‌های بی‌تجربه، بعضی اوقات در دور ریختن متغیرهایی که به نظرشان بدردنخور هستند و نگه داشتن تعداد کمی از متغیرهایی که به نظرشان مهم می‌رسند، بسیار شتاب دارند. روش‌های

داده‌کاوی این اجازه را می‌دهند که داده خودش مشخص کند که چه چیز مهم و چیز بی‌اهمیت است. اغلب، متغیرهایی که قبلاً از آنها چشم پوشی کرده‌ایم و آنها را کنار گذاشته‌ایم، در صورت ترکیب با متغیرهای دیگر، ارزش پیشگویانه پیدا می‌کنند. این درست است که مدل نهایی فقط براساس تعداد کمی متغیر بنا نهاده می‌شود. ولی این متغیرهای اندک گاهی از ترکیب چندین متغیر دیگر بدست می‌آیند و ممکن است که در ابتدا قابل مشاهده نباشند که یکی از آنها در نهایت مهم خواهد بود.

داده باید حاوی چه چیزی باشد ؟

داده در کمترین حالت باید حاوی تمام مثال‌های ممکن از نتایج مطلوب باشد. در داده‌کاوی مستقیم، که هدف پیش‌بینی متغیر هدف است، داشتن یک مجموعه مدل دربردارنده داده‌های دسته‌بندی شده، بسیار حساس است. برای تمایز بین افرادی که دوست دارند وام دریافت کنند و افرادی که دوست ندارند، از هر دو دسته مثال‌های متفاوت زیادی لازم است تا مدلی ساخته شود که یکی را از دیگری تشخیص دهد. وقتی یک درخواست جدید می‌رسد، درخواست او با مشتریان گذشته مقایسه می‌شود، یا بطور مستقیم در استدلال مبتنی بر حافظه، یا بطور غیرمستقیم از طریق قوانین یا شبکه‌های عصبی برگرفته از داده‌های سابقه‌ای. اگر درخواست جدید "شبهه به نظر رسید" با درخواست‌هایی که در گذشته قرارداد شده‌اند، رد خواهد شد.

مرحله 3: پیش به سوی شناخت داده :

داده‌کاوان خوب، بر روی درک مستقیم خیلی تکیه می‌کنند. تنها راه برای توسعه درک یک مجموعه داده ناآشنا، غوطه‌ور شدن در آن است. از این راه شما بسیاری از مسائل کیفیت داده را کشف خواهید کرد و سوالات بسیاری برایتان مطرح می‌شود.

آزمودن توزیع‌ها :

به هر چیزی که شما را متحیر می‌کند، توجه کنید؛ اگر یک متغیر کد ناحیه داریم، آیا "کالیفرنیا"، بیشترین طول را دارد؟ اگر متغیر جنسیت داریم، آیا تعداد زنان و مردان یکی است؟ به دامنه هر

متغیری دقت کنید. ابزارهای بصری‌سازی، می‌تواند در راه کنکاش اولیه پایگاه‌داده، کمک‌کننده باشد.

مقایسه مقادیر با توصیف‌ها :

به مقادیر هر یک متغیرها دقت کنید و مقایسه کنید که آیا با توصیف آن متغیر در مستندات، مطابقت دارد؟ این تمرین مشخص می‌کند که آیا مستندات درست و کامل هستند؟

اعتبارسنجی فرضیات :

با استفاده از ابزارهای بصری‌سازی، مثل گراف‌های خطی، نگاشت‌ها و نقشه‌های پراکنده، می‌توان فرضیاتی را درباره داده‌ها، بررسی و اعتبارسنجی کرد. به رابطه متغیر مقصد و سایر متغیرها دقت کنید تا چیزهایی از این قبیل را ببینید، مثل رابطه بین درآمد و جنسیت.

پرسیدن سوالات بسیار:

یک خروجی مهم از روال جستجوی داده‌ها، فهرستی از سوالاتی است از افرادی که داده‌ها را تهیه می‌کنند. اغلب این سوالات، بررسی‌های بیشتری را می‌طلبد، چون کاربران کمی، به همان دقت داده‌کاو، داده را بررسی می‌کنند.

مرحله 4: ساختن یک مجموعه مدل :

مجموعه مدل، شامل تمام داده‌هایی است که در فرایند مدلسازی بکار می‌روند. یکسری از این داده‌ها، برای یافتن الگوها، برخی دیگر برای اعتبارسنجی مدل، و برخی دیگر برای ارزیابی کارایی مدل، بکار می‌روند. ساختن یک مجموعه مدل، نیاز به اسمبل کردن داده از منابع چندگانه و سپس آماده کردن آنها برای تحلیل، دارد.

گردآوری امضاهای مشتری :

مجموعه مدل، یک جدول یا مجموعه‌ای از جداول با یک ردیف برای هر آیتم و فیلدهایی برای هر آنچه درباره آن آیتم می‌دانیم، است. وقتی داده مشتریان را توصیف می‌کند، ردیف‌های

مجموعه مدل اغلب "امضاهای مشتری" نامیده می‌شوند. گردآوری این اطلاعات اغلب نیاز به پرسش‌های پیچیده از چندین جدول و سپس یکپارچه‌سازی آنها با داده‌های دیگر منابع، دارد.

ساختن یک نمونه متعادل :

اغلب وظیفه داده‌کاوی شامل یادگیری تمایز بین گروه‌های مختلف است: مثل پاسخگوها و پاسخ‌نندگان، خوب‌ها و بد‌ها، و غیره. اگر تعداد اعضای این گروه‌ها متناسب باشد، الگوریتم‌های داده‌کاوی بهتر عمل می‌کنند. بنابراین قبل از مدلسازی، مجموعه داده باید با نمونه‌سازی از گروه‌های مختلف با نرخ‌های متفاوت و افزودن یک فاکتور وزن دهی، متوازن شود.

شامل چارچوب‌های زمانی چندگانه باشد :

هدف اصلی یک متدلوژی، ساختن یک مدل پایدار است. یکی از معانی آن این است که مدل‌ها برای هر زمانی از سال کار کنند و برای آینده مناسب باشند. این اتفاق زمانی می‌افتد که در مجموعه مدل، تمام داده‌ها از یک دوره خاص زمانی از سال، نیامده باشند. ساختن مدل براساس یک دوره زمانی خاص، این خطر را دارد که الگوهایی که عموماً درست نیستند را آموزش دهد. یک مثال جالب برای این مسئله این است که یکبار یک قانون وابستگی را برای داده‌های یک هفته سوپرمارکتی بدست آوردند که تمام قانون‌ها به تخم مرغ، ختم می‌شد. یعنی پیش‌بینی سبد خرید، بیان می‌کرد که هر کسی چیزی می‌خرد، تخم مرغ هم می‌خرد! البته این نباید زیاد تعجب برانگیز باشد اگر بدانیم که مجموعه مدل ما مربوط به هفته قبل از عید پاک بوده است!

ساختن یک مجموعه مدل برای پیش‌بینی :

اگرچه ممکن است که مجموعه مدل شامل فریم‌های مختلف زمانی باشد، اما هر امضای مشتری، یک فاصله زمانی بین متغیرهای پیش‌بینی کننده و متغیرهای هدف، خواهد داشت. همیشه می‌توان زمان را به سه دوره تقسیم کرد: گذشته، حال، آینده. وقتی یک پیش‌بینی انجام می‌شود، مدل، داده‌های گذشته را برای پیش‌بینی آینده بکار می‌گیرد. تمام این سه دوره زمانی باید در مجموعه مدل نشان داده شود. البته تمام داده‌ها در پایگاه داده‌ها از گذشته می‌آیند؛ بنابراین دوره‌های زمانی

در مجموعه مدل، بطور دقیق‌تر شامل گذشته دور، گذشته نچندان دور، و گذشته اخیر است. مدل‌های پیش‌بینی برای یافتن الگوهایی در گذشته دور ساخته شده‌اند که خروجی‌هایشان، گذشته اخیر را توصیف می‌کند. وقتی مدل مستقر شود، می‌تواند داده‌های گذشته اخیر را برای پیش‌بینی آینده بکار برد.

اگر یک مدل، داده‌های ماه June (گذشته نه چندان دور) را برای پیش‌بینی July (گذشته اخیر) بکار برد، نمی‌تواند تا زمانی که داده‌های آگوست در دسترس باشند، برای پیش‌بینی سپتامبر استفاده شود. ولی داده‌های آگوست کی در دسترس خواهد بود؟ قطعاً در خود آگوست نه. از آنجا که داده‌ها را باید جمع‌آوری، پالایش، بارگذاری و تست کرد، در بسیاری از شرکتها، داده‌های آگوست، حتی تا اوایل اکتبر در دسترس نیستند. راه‌حل شامل یک ماه تاخیر در مجموعه مدل است.

مرحله 5: تثبیت مسئله با داده‌ها :

تمام داده‌ها کثیف هستند. تمام داده‌ها مسئله دارند. آیا مسئله ای هست که با تکنیک‌های داده‌کاوی تغییر کند؟ برای مثال، برای درخت تصمیم‌گیری، مقادیر گم‌شده و چیزهای جدا، نمی‌تواند دردسر زیادی درست کند. برای بقیه، مثل شبکه‌های عصبی، می‌توانند باعث دردسرهای زیادی شوند.

متغیرهای طبقه‌بندی شده، با مقادیر بسیار :

متغیرهایی مثل کدپستی، کشور و کد شغل، همه از متغیرهایی هستند که اطلاعات مفیدی را بیان می‌کنند، ولی نه از راهی که بیشتر الگوریتم‌های داده‌کاوی می‌توانند مدیریت کنند. مقادیر ممکن بسیاری برای متغیرهایی که این اطلاعات را حمل می‌کنند وجود دارد، و مثال‌های کمی در داده‌های شما برای بیشتر این مقادیر وجود دارد.

متغیرهای این چینی، یا باید گروه‌بندی شوند که در آن صورت کلاس‌های بسیار زیادی برای رابطه آنها با متغیر پیش‌بینی، بوجود خواهد آمد؛ یا باید با خصوصیات مطلوب از کدپستی و

غیره جایگزین شوند. جایگزینی کدپستی با کدپستی میانگین، جایگزینی شغل با میانه حقوق برای شغل و به همین صورت.

توزیع مورب و پرت، مسائلی را برای هر تکنیک داده‌کاوی که مقادیر را بطور ریاضی بکار می‌برد، (مثلا با ضرب وزن آنها در یکدیگر و اضافه کردن آنها به هم) ایجاد می‌کند. در بسیاری از موارد، رکوردهایی که پرت افتاده اند، خارج می‌شوند. در سایر موارد، بهتر است که مقادیر را به اندازه دامنه یکسان تقسیم کرد. گاهی هم بهترین روش برای کاهش دامنه مقادیر این متغیرها، جایگزینی هر مقدار مثلا با لگاریتم آن است.

مقادیر گم شده :

برخی از الگوریتم‌ها، قادرند با مقادیر گم شده نیز مانند سایر مقادیر رفتار کنند و آنها را در قوانین بیاورند. متأسفانه برخی دیگر نمی‌توانند. دور ریختن تمام این رکوردهای با مقادیر گم‌شده، تبعیض را معمول می‌کند چون بعید است که این قبیل رکوردها بطور تصادفی توزیع شده باشند. جایگزینی آنها با مقادیری مثل میانه‌ها و مقدار متداول‌تر، اطلاعات جعلی را می‌سازد. جایگزینی آنها با مقادیر نامطلوب، حتی می‌تواند بدتر باشد، چون الگوریتم‌های داده‌کاوی آنها را نمی‌توانند تشخیص دهند، مثلا جایگزینی مقدار 999 برای متغیر سن. وقتی مقادیر گم‌شده باید جایگزین شوند، بهترین راه‌حل این است که با ساختن یک مدل که مقادیر گم‌شده بعنوان متغیر هدف باشند، آنها را مستند کنیم.

مقادیری که مفاهیم آنها در طول زمان تغییر کرده است :

وقتی داده‌ها از چندین نقطه از گذشته، می‌آیند، خیلی متداول است که مفهوم برخی مقادیر برای برخی فیلدها، در طول زمان تغییر کرده باشند. ممکن است که کلاس A اعتبار، همیشه بهترین کلاس باشد، اما آن دامنه رتبه‌های اعتباری که کلاس A را تشکیل می‌دهند، از زمانی به زمان دیگر، ممکن است تغییر کند. بحث کردن درباره این، کاملاً به یک انبار داده خوب طراحی شده نیاز دارد که این تغییرات در معنی بصورت یک متغیر جدیدی که می‌تواند در طول زمان معنای ثابتی داشته باشد، ثبت شده‌اند.

کدگذاری ناسازگار داده :

وقتی اطلاعات با موضوع یکسان از چندین منبع گردآوری شوند، اغلب ممکن است منابع مختلف روش های ارائه مختلفی برای داده‌ها داشته باشند. اگر این تفاوت‌ها دریافت نشوند، تمایزهای جعلی را اضافه می‌کنند که منجر به استنتاجات نادرستی می‌شود. با دانستن این تفاوت‌ها لازم است که داده‌ها به شکل رایج آنها ثبت شوند.

مرحله 6: تبدیل داده برای آوردن اطلاعات به سطح :

وقتی که داده‌ها همگذاری شدند و مسائل تثبیت شدند، هنوز داده باید برای تحلیل آماده شود. این شامل افزودن فیلدهای مشتق شده برای آوردن اطلاعات در سطح است. همچنین ممکن است این کار حذف پرت‌افتاده‌ها، گروه بندی کلاس‌ها برای رده بندی متغیرها، اعمال تبدیلاتی از قبیل لگاریتم‌ها و سایر کارهای مشابه باشد. آمایش داده‌ها، مسئله مهمی است که در کتابی بعنوان آمایش داده‌ها برای داده‌کاوی، بطور مفصل شرح داده شده است.

گرفتن گرایشات :

بیشتر داده‌های یکی شده، شامل سری‌های زمانی هستند. اطلاعات صورتحساب‌های عکس‌های فوری ماهانه، کارکردها، تماس‌ها و غیره. بیشتر الگوریتم‌های داده‌کاوی، سری‌های زمانی را درک نمی‌کنند. سیگنال‌هایی مثل "سه ماه کاهش سود سهام" نمی‌توانند نقطه به نقطه با مشاهدات هر ماه بطور مستقل، مطابقت داشته باشد. این تاحدودی به داده‌کاو بستگی دارد که اطلاعات رفتاری را با افزودن متغیرهای مشتق شده‌ای، به سطح بیاورد، مثل نسبت صرف شده در ماه‌های اخیر به ماه‌های قبلی، برای داشتن یک گرایش کوتاه مدت، و نسبت ماه‌های اخیر به همان ماه‌ها در سال قبل برای بدست آوردن یک گرایش بلندمدت.

ساختن نسبت‌ها و سایر ترکیبات متغیرها :

گرایشات یکی از مثال‌های آوردن اطلاعات در سطح با ترکیب متغیرهای چندگانه است. از اینها تعداد زیادی وجود دارد. اغلب این فیلدهای اضافی از روشهای موجود از طریق تحلیل‌گران دارای دانش گرفته شده‌اند، ولی کمتر توسط نرم‌افزارها توجه شده‌اند.

فیلدهای اضافه شده که روابط مهم از نظر خبره‌های آن زمینه را نشان می‌دهد، راهی است که اجازه می‌دهد که روال کنکاش از مهارتها بهره گیرد.

تبدیل شمارش‌ها به خصوصیات :

بسیاری از مجموعه داده‌ها، شامل مقادیر شمارش یا دلار هستند که خودشان بویژه جالب نیستند چون مطابق با سایر مقادیر، تغییر می‌کنند. خانواده‌های بزرگتر به نسبت خانواده‌های کوچکتر، پول بیشتری را در بقالی‌ها صرف می‌کنند. پول بیشتری در گوشت، محصولات، در بسته بندی بهتر، در محصولات شوینده، و همه چیز. پس مقایسه فیلد قیمت دلاری صرف شده بوسیله خانواده‌های با جمعیت متفاوت، در یک طبقه، مثل نانوايي، فقط آشکار می‌کند که خانواده‌های بزرگتر، بیشتر مصرف می‌کنند. مطلوبتر آن است که درجه مصرف هر خانواده را در هر طبقه، بررسی کنیم.

مرحله 7: ساختن مدلها :

جزئیات این مرحله از تکنیکی به تکنیک دیگر تغییر می‌کند. به عبارت کلی، این مرحله‌ای است که بیشتر کارهای ساخت یک مدل در آن اتفاق می‌افتد. در داده‌کاوی جهتدار، مجموعه آموزشی برای تولید تفسیری از متغیرهای مستقل یا هدف بر اساس متغیرهای وابسته یا ورودی، استفاده می‌شود. این تفسیر ممکن است یک شکل شبکه عصبی را بگیرد، یا درخت تصمیم‌گیری، یا یک گراف پیوند، یا سایر ارائه‌های دیگر از رابطه بین هدف و سایر فیلدهای پایگاه داده. در داده‌کاوی غیرجهتدار، هیچ متغیر هدفی وجود ندارد. مدل، روابط بین رکوردها را می‌یابد و آنها را بصورت قوانین وابستگی یا با تخصیص آنها به خوشه‌های متداول، بیان می‌کند.

ساختن مدل‌ها، اولین گام از روال داده‌کاوی است که بدرستی با نرم‌افزارهای داده‌کاوی مدرن بصورت اتوماتیک درآمده است. به همان دلیل، زمان کمتری را در پروژه داده‌کاوی، می‌گیرد.

مرحله 8: ارزیابی مدل ها :

این مرحله تعیین می کند که مدل ها کار می کنند یا نه. ارزیابی مدل ها باید به سوالات از این قبیل، پاسخ دهد:

- 1- چقدر یک مدل دقیق و قابل درک است؟
 - 2- مدل چقدر توانسته داده های مشاهده شده را تشریح کند؟
 - 3- چه میزان اطمینان از پیش بینی های مدل وجود دارد؟
- البته پاسخ به این سوالات به نوع مدل ساخته شده بستگی دارد.

ارزیابی مدل های توصیفی :

قانون، "اگر (ایالت = 'MA') آنگاه منبع گرما نفت است"، از قانون، "اگر (ناحیه = 339 یا 351 یا 413 یا 508 یا 617) باشد آنگاه منبع گرمایش نفت خواهد بود"، بسیار توصیفی تر است. حتی اگر دو قانون معادل هم باشند، اولی به نظر می رسد که توصیفی تر است.

قدرت بیان ممکن است بطور ذهنی به نظر برسد، ولی در اینجا در حقیقت یک راه تئوری برای اندازه گیری آن وجود دارد، که کوچکترین طول توصیف (MDL) نامیده می شود. برای یک مدل، MDL، عبارت است از تعداد بیت هایی که آن مدل برای کدگذاری قانون و تمام استثناءهای آن، می گیرد. قانونی بهتر است که به بیت های کمتری نیاز دارد. برخی ابزارهای داده کاوی MDL را استفاده می کنند تا ببینند که کدام مجموعه قانون را نگه دارند و کدام را دور بریزند.

ارزیابی مدل های جهتدار :

مدل های جهتدار در دقتشان در داده های نادیده اولیه، سنجیده می شوند. ارزیابی هر مدلی به زمینه آن بستگی دارد. مدل های یکسان می توانند برطبق یک پیمان، بد و در دیگری خوب باشند. محققان یک هدف از ساخت مدل ها دارند و آن امکان درک تمامیت مسئله است.

ارزیابی مدل می تواند در سطح کل مدل باشد یا در سطح پیشگوهای تکی. دو مدل با دقت های سراسری یکسان، ممکن است سطوح مختلفی از ناسازگاری را در پیشگوهای تکی شان داشته

باشند. یک درخت تصمیم‌گیری، برای نمونه، یک نرخ خطای دسته‌بندی سراسری دارد، ولی هر شاخه و برگ درخت نرخ خطای خودش را دارد.

ارزیابی طبقه‌بندها و پیشگوها :

برای وظایف پیشگویی و دسته‌بندی، دقت با عبارت نرخ خطا، سنجیده می‌شود. نرخ خطای دسته‌بندی روی یک مجموعه آزمایشی از پیش دسته‌بندی شده، بکار می‌رود برای تخمین نرخ خطا هنگامی که رکوردهای جدید دسته‌بندی می‌شوند. البته این روال فقط هنگامی معتبر است که مجموعه آزمایشی در نسبت بزرگتری ارائه شوند.

ارزیابی تخمین‌گرها :

برای وظایف تخمین، دقت می‌تواند با عبارت تفاوت بین نمره پیش‌بینی شده و اندازه واقعی، بیان شود. دقت هر جزئی که تخمین زده می‌شود و دقت کل مدل، هر دو مطلوب و مورد نظر ماست. ممکن است مدلی برای یک دامنه از مقادیر ورودی، کاملاً صحیح باشد و برای دامنه دیگری نباشد. روش استاندارد برای توصیف دقت مدل تخمین، اندازه‌گیری میزان دوری تخمین از میانگین. ولی می‌توان بسادگی با کم کردن مقدار تخمین از مقدار واقعی یا با میانه‌گیری برای مقادیر بی‌معنی، آنرا بدست آورد. "واریانس کوچکتر بمعنی دقت بالاتر تخمین است."

مقایسه مدل‌ها با استفاده از lift :

مدل‌های جهتدار، که با شبکه‌های عصبی، الگوریتم ژنتیک و درخت تصمیم همگی برای به انجام رساندن برخی وظایف ساخته شده‌اند. چرا نباید آنها درباره توانایی هایشان برای دسته‌بندی، تخمین و پیش‌بینی، مقایسه شوند؟ یک روش بسیار معمول برای مقایسه کارایی دسته‌بندی مدل‌ها، استفاده از نسبتی بنام Lift است. این مقیاس می‌تواند برای مقایسه مدل‌های طراحی شده برای سایر وظایف نیز بخوبی تطبیق داده شود. درحقیقت lift تغییر تمرکز از یک کلاس خاص را وقتی

مدل برای انتخاب یک گروه از جمعیت کل، استفاده می‌شود، اندازه می‌گیرد.

$$\text{lift} = P(\text{class}_i | \text{sample}) / P(\text{class}_i | \text{population})$$

فرض کنید می خواهیم مدلی را بسازیم برای پیش بینی اینکه چه کسی تمایل به پاسخ به اشتراک پست مستقیم دارد. معمولاً، ما مدلی می سازیم با استفاده از یک مجموعه داده آموزشی پیش طبقه بندی شده، و اگر لازم شد، یک مجموعه اعتبارسنجی پیش طبقه بندی شده. حالا ما آماده ایم که مجموعه تست را برای محاسبه lift مدل بکار ببریم.

دسته بندی کننده، رکوردها را در مجموعه تست بعنوان "پیش بینی شده برای پاسخ" یا "پیش بینی نشده برای پاسخ" دسته بندی می کند. البته این همیشه درست نیست، ولی اگر مدل به هیچ وجه مزیتی نداشته باشد، گروه علامت گذاری شده بعنوان "پیش بینی شده برای پاسخ"، شامل یک نسبت بالاتر از پاسخ دهندگان واقعی است به نسبت کل مجموعه تست. به این رکوردها توجه کنید. اگر مجموعه تست شامل 5 درصد از پاسخ دهندگان واقعی باشد و نمونه، شامل 50 درصد از پاسخ دهندگان واقعی باشد، lift مدل 10 است. (50/5)

آیا مدلی که Lift بیشتری داشته باشد، بهتر است؟ البته فهرستی از افراد، که نیمی از کل آنها پاسخ دهنده خواهند بود، به فهرستی که فقط یک چهارم آنها پاسخ دهنده اند، ترجیح داده می شود. درست؟ نه لزوماً - ممکن است که آن لیست فقط شامل 10 نام از پاسخ دهندگان باشد.

نکته اینجاست که lift یک تابع از اندازه نمونه است. اگر یک دسته بند فقط 10 تا از پاسخ دهندگان را بردارد و این 100 درصد زمان باشد، به مقدار $lift = 20$ خواهد رسید، بیشترین مقدار lift ممکن وقتی که جمعیت حاوی 5 درصد پاسخ دهندگان باشد.

مرحله 9: استقرار مدل ها :

استقرار مدل به معنی بردن آن از محیط داده کاوی به محیط رتبه بندی است. فرایند ممکن است ساده یا سخت باشد. در بدترین حالت (که در بیش از یک شرکت دیده شده) مدل در یک محیط مدلسازی خاص با استفاده از نرم افزاری که هیچ جایی اجرا نمی شود، توسعه داده می شود.

یک مسئله شایع دیگر این است که مدل از متغیرهای ورودی ای استفاده کرده باشد که در داده واقعی وجود ندارد. این مسئله نباید وجود داشته باشد از آنجا که متغیرهای ورودی مدل از فیلدهای برگرفته از داده اصلی در مجموعه مدل، گرفته شده اند.

مرحله 10: ارزیابی نتایج

مرحله 11: شروع دوباره :

از هر پروژه داده‌کاوای پرسش‌هایی بسیاری برمی‌خیزد. این چیز بسیار خوبی است. این به این معنی است که روابط جدیدی اکنون دیده می‌شوند که قبلاً دیده نمی‌شدند. روابط جدید کشف شده، فرضیه‌های جدیدی را پیشنهاد می‌کنند که باید آزمایش شوند و بنابراین روال داده‌کاوای از ابتدا آغاز می‌شود.

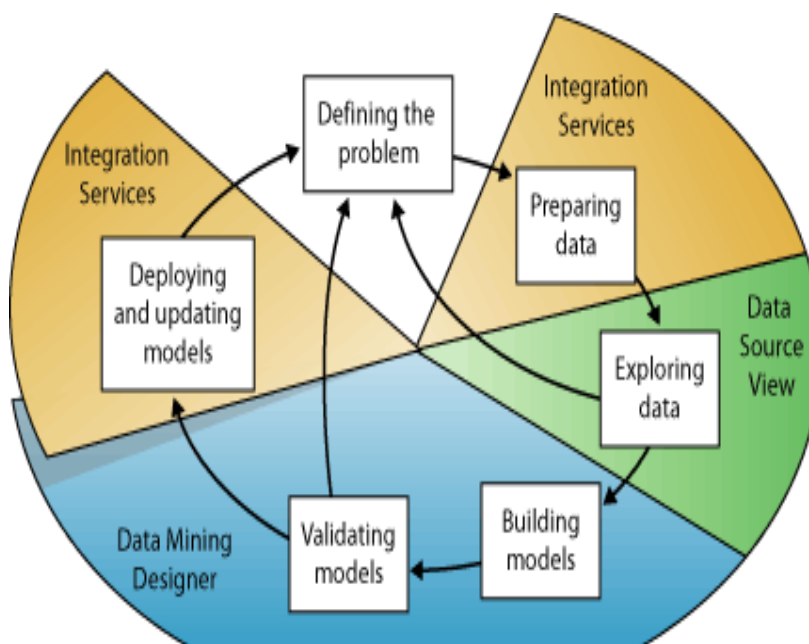
۲-۴-۱ دیگر نظریه‌ها (مراحل داده‌کاوای) :

- 1- تعیین اطلاعات گذشته
- 2- تمیز دادن و پردازش‌های اولیه. در این مرحله خطاهای داده تصحیح می‌شوند و داده‌های اشتباه جایگزین می‌شوند. این مرحله ممکن است تا 60 درصد از زمان داده‌کاوای را به خود اختصاص بدهد
- 3- یکپارچه‌سازی داده‌ها. معمولاً داده‌ها از منابع متفاوتی جمع‌آوری می‌شوند و باید به صورتی در آیند که یک مخزن داده‌ی مناسب را ایجاد کنند تا بتوان عملیات داده‌کاوای را بهتر انجام داد
- 4- انتخاب مجموعه‌ای از داده‌های هدف
- 5- یافتن ویژگی‌های مورد استفاده
- 6- نمایش داده‌ها به صورتی که بتوان برای داده‌کاوای مورد استفاده قرار گیرند
- 7- انتخاب عمیقات داده‌کاوای (طبقه‌بندی شده، خوشه‌بندی و ...)
- 8- انتخاب روش داده‌کاوای (شبکه‌های عصبی، درخت تصمیم و ...)
- 9- داده‌کاوای برای یافتن الگوی مناسب
- 10- ارزیابی و تحلیل الگوی بدست آمده و حذف الگوی نامناسب
- 11- تفسیر نتایج داده‌ها و استنتاج اطلاعات با ارزش

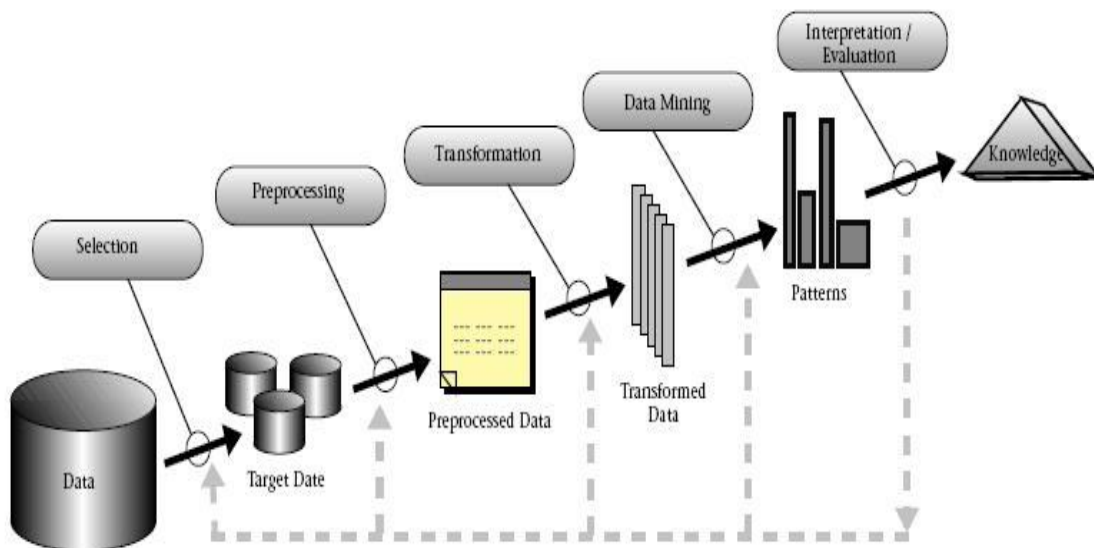
باید توجه داشت که جمع‌آوری و حفاظت از داده‌ها نکته بسیار مهمی است. اصولاً چون قالب و نوع داده‌ها در طول زمان تغییر می‌کند ممکن است بسیاری از داده‌های موجود در قالب‌های متفاوتی باشد و هم‌چنین به بسیاری از داده‌های قدیمی از بین رفته و دور ریخته می‌شود در حالی که ممکن است اهمیت این داده‌ها از داده‌های جدید به هیچ وجه کمتر نباشد. هم‌چنین

به علت اینکه داده ها میتوانند از منابع متلفی داخلی و خارجی مانند: کارکنان شرکت، مدیران شرکت، مشتریان، کار فرمایان و کارکنان باشند بازهم ممکن است قالب داده ها با هم یکسان نباشند. به همین دلیل انتخاب داده های درست و یکپارچه سازی قالب آن ها به منظور استفاده در داده کاوی از اهمیت بالایی برخوردار است. در شکل 2 می توان مراحل داده کاوی را به اختصار مشاهده کرد

- 1- درک قلمرو یا بیان مسئله و فرموله کردن فرضیه
- 2- انتخاب و جمع آوری داده ها
- 3- تبدیل داده ها
- 4- کاوش در داده ها
- 5- تفسیر نتیجه یا تفسیر مدل و رسیدن به نتایج



مراحل داده کاوی عموم-شکل 1



مراحل داده کاوی آکادمیک - شکل 2

۲-۵ وظایف داده کاوی. Error! Reference source not found. :

داده کاوی می تواند برای حل صدها مسئله کسب و کار بکار آید. بر اساس طبیعت این مسائل، ما آنها را به وظایف داده کاوی زیر دسته بندی می کنیم.

1- دسته بندی :

دسته بندی یکی از رایج ترین وظایف داده کاوی است. مسائل تجاری ای مانند تحلیل رویگردانی، مدیریت ریسک و هدف گیری موردی، شامل دسته بندی می شوند.

دسته بندی به تخصیص موارد به گروه ها براساس یک خصوصیت قابل پیش بینی، برمی گردد. هر مورد شامل مجموعه ای از خصوصیات است که یکی از آنها خصوصیت کلاس (خصوصیت قابل پیش بینی) نام دارد. این وظیفه شامل یافتن مدلی است که خصوصیت کلاس را براساس تابعی از خصوصیات ورودی تشریح کند. برای آموزش مدل رده بندی، شما به مقدار کلاس موارد ورودی در مجموعه داده های آموزشی نیاز دارید. الگوریتم های داده کاوی که به یک هدف برای یادگیری نیاز دارند، الگوریتم های نظارت شده، نام دارند.

الگوریتم‌های رایج دسته‌بندی شامل درخت‌های تصمیم، شبکه عصبی، و ناوی بیز است.

2- خوشه‌بندی :

خوشه‌بندی که گاهی قطعه‌بندی نیز خوانده می‌شود برای شناسایی گروه‌بندی طبیعی موارد بر اساس یک مجموعه خصوصیت بکار می‌رود. موارد درون یک گروه بیشترین یا کمترین میزان شباهت را در مقادیر خصوصیات دارند.

الگوریتم خوشه‌بندی مجموعه داده را براساس این دو خصوصیت به سه گروه تقسیم کرده است.

خوشه‌بندی یک وظیفه داده‌کاوی بدون نظارت است. هیچ خصوصیت تکی برای هدایت روال آموزش بکار نمی‌رود و تمام خصوصیات ورودی بطور یکسان هدف هستند. بیشتر الگوریتم‌های خوشه‌بندی، مدل را در تعداد تکرارهای معینی می‌سازند و هنگامی که مدل همگرا شد، متوقف می‌شوند و آن هنگامی است که مرزهای گروه‌ها تثبیت شد.

3- تخمین:

تخمین با نتایجی که با ارقام پیوسته نشان داده شده‌اند سروکار دارد. در تخمین داده‌های ورودی داده می‌شوند و به رقمی در متغیرهای ناشناس مداوم چون درآمد یا تراز کارت اعتباری ختم می‌شود. در عمل، تخمین اغلب برای انجام دسته‌بندی استفاده می‌شود.

یک شرکت کارتهای اعتباری که مایل است یک فضای تبلیغاتی را در پاکتهای صورتحساب به یک تولید کننده پوتین اسکی بفروشد، باید مدل دسته‌بندی تهیه کند که همه دارندگان کارتها را در یک یا دو دسته قرار دهد: اسکی باز یا غیر اسکی باز.

روش دیگر ایجاد مدلی است که به هر دارنده کارت، یک امتیاز تمایل به اسکی تخصیص می‌دهد. این ارقام می‌تواند صفر و یک باشد که نشانگر احتمال تخمین زده شده است که دارنده کارت یک اسکی باز است یا نه. عمل دسته‌بندی اکنون به ایجاد امتیازی آستانه‌ای منجر می‌گردد. هر کسی که امتیازی بیشتر یا مساوی با امتیاز آستانه داشته باشد به عنوان اسکی باز قلمداد می‌شود. هر کسی که امتیازی کمتر از امتیاز مورد نظر را داشته باشد اسکی باز محسوب نمی‌شود.

گردد. شرکت تولید پوتین‌های اسکی برای پست پانصد هزار اوراق تبلیغاتی بودجه‌ریزی نموده است. اگر از روش دسته‌بندی استفاده شده و یک و نیم میلیون نفر اسکی باز تعیین شده‌اند پس به راحتی می‌توان به صورت تصادفی تبلیغات را در صورتحسابهای پانصد هزار نفر منتخب از آن افراد قرارداد. اگر از طرف دیگر هر دارنده کارت، امتیاز تمایل به اسکی را داشته باشد می‌توان تبلیغات را برای پانصد هزار از محتمل‌ترین کاندیداها فرستاد.

مثالهایی دیگر از تخمین:

- 1- تخمین تعداد فرزندان در یک خانواده
- 2- تخمین درآمد کل یک خانواده
- 3- تخمین عمر یک مشتری
- 4- تخمین احتمال اینکه فردی به یک تقاضای تغییر میزان اعتبار بانکی پاسخ دهد.

4- وابستگی:

وابستگی یکی دیگر از وظایف معمول داده‌کاوی است که تحلیل سبد بازار نیز خوانده می‌شود. یک مسئله کسب و کار معمول وابستگی، تحلیل یک جدول تراکنش فروش و شناسایی آن محصولاتی است که اغلب در یک سبد خرید یکسان قرار می‌گیرند. استفاده رایج از وابستگی شناسایی مجموعه رایج اقلام خرید (مجموعه اقلام متناوب) و قوانینی برای فروش-مقاطع است.

در عبارات وابستگی، هر محصولی، یا بطور کلی هر جفت خصوصیت/مقداری، یک قلم در نظر گرفته می‌شود. وظیفه وابستگی دو هدف دارد: یافتن مجموعه اقلام متناوب و یافتن قوانین وابستگی.

بیشتر الگوریتم‌های وابستگی برای یافتن مجموعه اقلام متناوب، مجموعه داده‌ها را چندین بار پیمایش می‌کنند. آستانه تناوب (پشتیبانی)، قبل از شروع پردازش، توسط کاربر تعیین می‌شود. برای مثال پشتیبانی برابر با 2 درصد به معنی است که مدل فقط آن اقلامی که در 2 درصد از کارت‌های خرید دیده می‌شوند را تحلیل کند. یک مجموعه قلم می‌تواند مانند مجموعه {محصول="پپسی"، محصول="چیپس"، محصول="آبمیوه"} دیده شود. هر مجموعه قلم یک اندازه دارد که تعداد آیت‌های آن است. مثلاً اندازه این مجموعه آیت‌ها برابر با 3 است.

جدای از شناسایی مجموعه اقلام بر مبنای پشتیبانی، بیشتر الگوریتم‌های وابستگی، قوانین را نیز می‌یابند. یک قانون وابستگی شکلی مانند $A, B \Rightarrow C$ با یک احتمال دارد، که A و B و C همگی مجموعه اقلام هستند. همچنین احتمال در ادبیات داده‌کاوی اعتماد نیز خوانده می‌شود. احتمال یک مقدار آستانه است که یک کاربر قبل از آموزش یک مدل وابستگی، نیاز دارد که تعیین کند. برای مثال، این یک قانون معمول است: محصول = "پسی" و محصول = "چیپس" = <= محصول = "آمیوه" با احتمال 80%. تفسیر این قانون بسیار سراسر است. اگر یک مشتری پسی و چیپس خرید، با احتمال 80 درصد ممکن است که او آمیوه نیز بخرد.

5- رگرسیون :

وظیفه رگرسیون مشابه با دسته‌بندی است. تفاوت اساسی در خصوصیت پیش‌بینی است که یک عدد پیوسته است. تکنیک رگرسیون سالها در حوزه آمار مطالعه شده است. رگرسیون خطی و منطقی از روش‌های بسیار رایج رگرسیون هستند. سایر تکنیک‌های رگرسیون شامل درخت‌های رگرسیون و شبکه‌های عصبی است.

وظیفه رگرسیون می‌تواند بسیاری از مسائل کسب و کار را حل کند. برای مثال، آنها می‌توانند در پیش‌بینی نرخ‌های خریداری و آزادسازی کوپن بر اساس ارزش وجه، روش توزیع و حجم توزیع یا پیش‌بینی سرعت باد بر اساس دما، فشار هوا و رطوبت، بکار روند.

6- پیشگویی :

پیشگویی یک وظیفه داده‌کاوی مهم دیگری است. ارزش سهام MSFT فردا چگونه خواهد بود؟ مقدار فروش پسی در ماه آینده چگونه خواهد بود؟ پیشگویی می‌تواند به این سوالات پاسخ دهد.

در عمل پیش‌بینی، تنها روش برای بررسی صحت مدل، دیدن آینده و مقایسه نتیجه مدل و پدیده واقع شده می‌باشد. هر یک از تکنیک‌های استفاده شده در دسته بندی و تخمین را می‌توان برای استفاده در پیش‌بینی تطبیق داد. از داده‌های پیشین برای تهیه یک مدل که بیانگر رفتار مشاهده

کنونی است استفاده می‌شود. وقتی این مدل برای ورودی های کنونی به کار رفت نتیجه کار پیش‌بینی رفتار آینده خواهد بود.

7- تحلیل توالی :

تحلیل توالی برای یافتن الگوها در رشته‌های گسسته بکار می‌رود. یک توالی از مقادیر (یا نواحی) گسسته ترکیب شده است. برای مثال، یک توالی DNA یک رشته مرکب از چهار ناحیه مختلف A، G، C و T است. یک توالی کلیک وب شامل رشته‌هایی از URLهاست. خریدهای کامپیوتری را هم می‌توان بصورت داده‌های توالی مدل کرد. برای مثال، یک مشتری ابتدا یک کامپیوتر، سپس اسپیکرها، و در نهایت یک دوربین وبی می‌خرد. تفاوت داده‌های توالی و رشته‌های زمانی در این است که رشته‌های زمانی شامل اعداد پیوسته و توالی شامل مقادیر گسسته است.

داده‌های توالی و وابستگی در اینکه هر یک شامل یک مجموعه آیت‌ها یا ناحیه هستند، شبیه می‌باشند. تفاوت آنها در این است که مدل‌های توالی انتقال‌های نواحی را تحلیل می‌کنند و مدل وابستگی فرض می‌کند که هر فقره در یک کارت خرید برابر یا مستقل باشد. با مدل توالی، خرید یک کامپیوتر قبل از اسپیکرها با خرید اسپیکرها قبل از کامپیوتر متفاوت است. درحالی که با یک الگوریتم وابستگی این دو توالی خرید، مجموعه اقلام یکسانی را می‌سازند.

تحلیل توالی، نسبتاً یک وظیفه داده‌کاوی جدید است. و در دو نوع کاربرد بسیار مهم شده است: تحلیل ثبتهای وب و تحلیل DNA. امروزه تکنیک‌های توالی مختلفی مانند زنجیره مارکوف در دسترس ما هستند. محققان در این زمینه بطور فعالی در حال تحقیق هستند

8- تحلیل انحراف:

تحلیل انحراف برای یافتن موارد نادری است که با دیگران بسیار متفاوت هستند. این که همچنین کشف outlier نیز نامیده می‌شود که تشخیص تغییرات معنی‌دار از رفتار مشاهده شده قبلی است. تحلیل انحراف می‌تواند در کاربردهای بسیاری استفاده شود. رایج‌ترین استفاده از آن در کشف

تقلبات کارتهای اعتباری است. شناسایی موارد غیرعادی از بین میلیون‌ها تراکنش یک چالش بسیار مشکل است. سایر کاربردها می‌تواند شامل کشف ورود سرزده به شبکه، تحلیل خطای تولید و غیره باشد.

هیچ استاندارد برای تحلیل انحراف وجود ندارد و این هنوز یک موضوع باز برای تحقیقات است. اما تحلیلگران معمولاً نسخه‌های اصلاح شده الگوریتم‌های درخت‌های تصمیم، خوشه‌بندی، یا شبکه‌های عصبی را برای این کار بکار می‌گیرند. بمنظور ساختن قوانین معنی‌دار، تحلیل‌گران نیاز به داشتن نمونه‌های بسیار از این موارد نادر در مجموعه‌های آموزشی خود دارند.

9- نمایه‌سازی:

گاهی اوقات هدف داده‌کاوی تنها توصیف آن چیزی است که در یک پایگاه داده پیچیده در جریان است. نتایج نمایه‌سازی درک ما را از مردم، محصولات یا فرآیندهایی که داده‌ها را در مرحله اول تولید کرده‌اند افزایش می‌دهد. توصیف خوب رفتار اغلب توضیح خوبی هم به همراه دارد.

شکاف جنسیتی مشهور در سیاست آمریکا مثالی از این دست است که چگونه این توصیف ساده که تعداد زنان حامی حزب دمکرات بیش از مردان است می‌تواند توجه بیشتر و مطالعات تکمیلی را برای روزنامه‌نگاران، جامعه‌شناسان، اقتصاددانان و دانشمندان علوم سیاسی ایجاد کند.

درخت‌های تصمیم ابزار مفیدی برای نمایه‌سازی می‌باشد. قوانین وابستگی و خوشه‌بندی را نیز می‌توان برای نمایه‌سازی‌ها استفاده نمود.

سه وظیفه نخست مثال‌هایی از داده‌کاوی جهت‌دار است. در داده‌کاوی جهت‌دار همیشه یک متغیر هدف وجود دارد برخی مواقع دسته‌بندی می‌شود، تخمین زده می‌شود یا پیش‌بینی می‌شود. روال ساختن یک دسته‌بند، با یک مجموعه پیش‌تعریف شده کلاس‌ها و مثالهایی از رکوردهایی شروع می‌شود که در حال حاضر بخوبی کلاس‌بندی شده‌اند. بطور مشابه، روال ساختن یک تخمین‌گر، با داده‌های سابقه‌ای شروع می‌شود که مقدار متغیر مقصد برای آنها در حال حاضر

مشخص است. وظیفه مدل سازی یافتن قوانینی است که مقادیر شناخته شده متغیر هدف را توضیح دهد.

در داده‌کاوی غیر جهتدار، هیچ متغیر هدفی وجود ندارد. گروه بندی وابستگی و خوشه بندی، داده‌کاوی بدون جهت هستند. وظیفه داده‌کاوی یافتن الگوهای موازی است که به هیچ تغییری مقید نیستند. یک شکل رایج داده‌کاوی غیرجهتدار، خوشه بندی است، که گروه‌هایی از رکوردهای مشابه را می‌یابد بدون هیچ دستورالعملی درباره اینکه به چه تغییری باید بیشتر توجه شود. داده‌کاوی غیر جهتدار با طبیعتش توصیف شده است.

بیان مسئله و فرموله کردن فرضیه :

در ابتدای امر پیش زمینه کشف دانش، فهم درست داده و مساله می باشد. بدون این فهم درست هیچ الگوریتمی صرف نظر از خبره بودن آن نمی تواند نتیجه مطمئنی برای شما حاصل نماید و داده را جهت کاوش آماده نموده یا نتایج را به طور صحیح تفسیر نمود. برای استفاده بهتر از داده کاوی باید یک بیان واضح از هدف داشت. در این مرحله آنچه نیاز است ترکیبی از تخصص یک زمینه کاربردی و یک مدل داده کاوی است و شاید بتوان گفت یک تقابل نزدیک سر یک مسئله واحد و چندین فرضیه فرموله شده بین متخصصین داده کاوی و متخصصین کاربردی میباشد.

1- انتخاب و جمع آوری داده ها:

این مرحله در ارتباط با چگونگی تولید و جمع آوری داده ها است.

بطور کلی، دو امکان وجود دارد:

روش آزمون طراحی: زمانی است که فرایند تولید داده ها تحت کنترل یک متخصص کاربردی (مدل ساز سیستم) باشد.

روش دیداری: امکان دوم زمانی مطرح است که متخصص قادر به تولید فرآیند نیست یعنی تولید داده بصورت تصادفی در نظر گرفته شود.

پس از اینکه داده‌ها جمع‌آوری شدند یا در فرایند جمع‌آوری داده‌ها تا اندازه‌ای قرار گرفتند، توزیع نمونه‌گیری کاملاً نامعلوم است. (یعنی داده‌هایی که بعداً برای تست و بکارگیری آن مدل بکار می‌روند از چند نمونه مشابه استفاده می‌شوند.)

نکته: برای فرایند داده‌کاوی داده‌های مورد نیاز موجود در انبار داده‌ها باید انتخاب شوند. درک این مطلب که برای ارزیابی یک مدل که بعداً برای تست و بکارگیری آن مدل بکار می‌رود، موفقیت‌آمیز باشد، بسیار مهم است در غیر اینصورت نتایج درستی حاصل نمی‌گردد.

مثلاً انبار داده‌ها شامل انواع مختلف و گوناگونی از داده‌ها است به عنوان مثال در یک پایگاه داده‌های مربوط به سیستم فروشگاهی، اطلاعاتی در مورد خرید مشتریان، خصوصیات آماری آنها، dispatcherها (توزیع‌کنندگان)، مشتریان، حسابداری و... وجود دارند که همه آنها در داده‌کاوی مورد نیاز نیستند.

2- پیش پردازش‌ها یا تبدیل داده‌ها

زمانی که داده‌های مورد نیاز از پایگاه داده‌های موجود در انبار داده‌ها "جمع‌آوری" شدند و داده‌های مورد کاوش مشخص گردیدند، معمولاً به تبدیلات خاصی روی داده‌ها نیاز است که شامل حداقل دو مرحله متداول می‌باشد:

آشکارسازی (حذف) داده‌های غیرعادی:

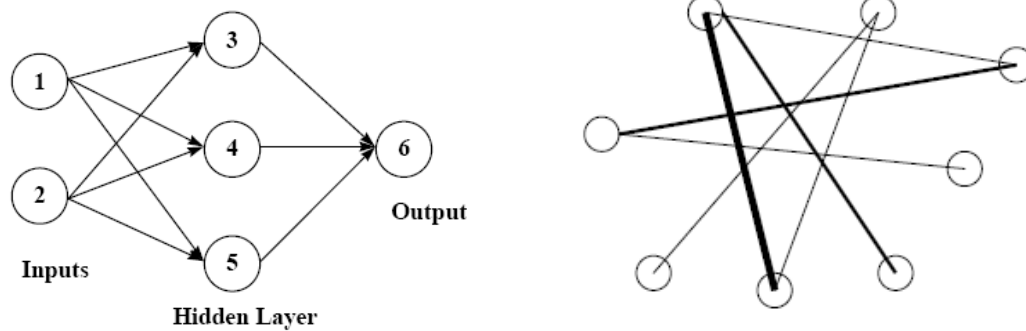
داده‌های غیرعادی یا غیر معمول در حقیقت داده‌های نتیجه‌سنجش خطاها، کدنویسی و ثبت خطاها است. در اینجا باید یا 1. داده‌های غیرعادی را تشخیص داد و حذف کرد یا 2. باید روش‌های قوی مدل‌سازی را بگونه‌ای توسعه داد که نسبت به این نوع داده‌ها غیر حساس باشند.

ویژگی‌های مقیاس‌بندی، رمزگذاری و انتخاب:

در تبدیل داده‌ها توصیه می‌شود که داده‌ها را جهت تحلیل و بررسی مقیاس‌بندی و رمزگذاری کرد. مثلاً یک مشخصه با دامنه [0 و 1] و دیگری با دامنه [1000 و -100] دارای ارزش مشابهی در تکنیک‌های اعلام شده نیستند. که در صورت نادیده گرفتن همین تفاوت در دامنه داده‌ها، روی نتایج نهایی داده‌کاوی تاثیر خواهند گذاشت.

3- برآورد مدل یا کاوش در داده ها

در این مرحله داده های تبدیل شده با استفاده از تکنیکها و عملیتهای داده کاوی مورد کاوش قرار می گیرند تا الگوهای مورد نظر کشف شوند. یا به عبارتی دیگه، انتخاب و پیاده سازی تکنیک های داده کاوی در این مرحله صورت میگیرد. البته این فرایند خیلی روشن و واضح نیست زیرا هنگام پیاده سازی ممکن است که مبتنی بر چندین مدل در یک فرآیند تکرار باشد. (این مدل ها بطور کامل تر در مباحث مربوط به مفاهیم انواع دسته بندی، درختان تصمیم و قوانین تصمیم، شبکه های عصبی، انواع الگوریتم ها و... پیاده سازی می شوند)



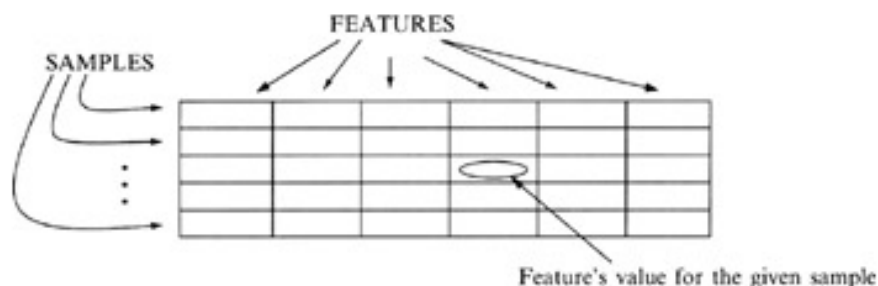
شکل 3 - یک شبکه عصبی

4- تفسیر نتیجه یا تفسیر مدل و رسیدن به نتایج

اطلاعات استخراج شده با توجه به هدف کاربر تجزیه و تحلیل شده و بهترین نتایج باید در تصمیم گیری کاربر موثر می باشند. هدف از این مرحله تنها ارائه نتیجه (بصورت منطقی و یا نموداری) نیست، بلکه پالایش اطلاعات ارائه شده به کاربر نیز از اهداف مهم این مرحله است.

۶-۲ آماده سازی داده ها

مدل استاندارد داده ها:



جدول داده خام - شکل 4

مدل استاندارد داده های ساخت یافته که حاوی فیلدهایی با مقادیر عددی می باشند برای داده کاوی دارای حالات مختلفی می باشند. در داده کاوی جهت نمایش اشیای ذخیره شده "خصیصه" یا "ویژگی" و جهت نمایش رکوردهای سطرها، از اصطلاح "نمونه" یا "حالت" استفاده می شود. در آماده سازی داده ها گاهی اوقات فقط به عنوان یک مرحله از فرایند داده کاوی تلقی می شود.

در آماده سازی داده ها دو وظیفه داریم:

1- ساماندهی داده ها به یک شکل استاندارد که برای پردازش بوسیله تکنیک های

داده کاوی آماده شوند (یک شکل استاندارد، یک جدول رابطه ای است).

2- آماده سازی مجموعه داده هایی که به بهترین عملکرد داده کاوی منجر شود.

تبدیل و تغییر وضعیت داده های خام

برای تبدیل یا تغییر وضعیت داده های خام ما باید از چند روش یا تکنیک تغییر داده ها (البته بسته به نوع داده ها) یکی را انتخاب کنیم.

* نرمال سازی

هدف از نرمالسازی حذف افزونگی داده و باقی نگهداشتن وابستگی بین داده های مرتبط از طریق ایجاد رابطه ستون های غیر کلیدی در هر جدول کلید است. این فرآیند اغلب باعث

ایجاد جداول بیشتر می شود ولی از این طریق اندازه پایگاه داده را کاهش داده و بهبود کارایی داده را تضمین می کند.

با توجه به وضعیت ممکن است داده ها از چند پایگاه داده نرمال شده استخراج شوند و در یک انبار داده غیر نرمال قرار گیرد. این روش برای مخزن داده Data warehouse استاندارد خوبی است.

BIG UGLY TABLE										
SaleNo	SaleDate	ProductNo	Qty	Amount	Salesrep	CustomerNo	First	Last	Address	CreditLimit
12345	Aug 12 2002	AQX88916	1	23.95	Dave Williams	4649-4673	Richard	Johnston	14 West Avenue	1000
12346	Aug 12 2002	AQX88916	7	167.65	Sara Thompson	1113-7741	Wayne	Jones	42 York Street	<null>
12347	Aug 13 2002	AHL46785	3705	5001.75	Li Qing	1166-3461	Amelia	Waverley	995 Forth Street	<null>
12348	Aug 13 2002	DHU69863	50	118.5	Sara Thompson	<null>	<null>	<null>	<null>	<null>
12349	Aug 14 2002	DHU69863	940	2227.8	Sara Thompson	1166-3461	Amelia	Waverley	995 Forth Street	<null>
12350	Aug 14 2002	DHU69863	42	99.54	Sara Thompson	7671-3496	Antonio	Gonzales	55B Granary Lane	<null>
12351	Aug 14 2002	AQX88916	55	1317.25	Dave Williams	6794-1674	Diane	Adams	364 East Road	1500

نمونه جدول استخراج داده- شکل 5

مقیاس دهی اعشار: مقیاس دهی اعشار نقطه اعشاری را انتقال می دهد اما بیشترین مقادیر اصلی را حفظ می کند. مقیاس کلی و انتخابی، مقادیر را در دامنه 1- تا 1 برقرار میکند.

$$V'(i) = V(i)/10^K$$

معادله فوق به ازای کوچکترین مقدار k با این فرض که $1 < |V'(i)|$ باشد، تعریف شده است.

نرمال سازی حداقل: حداکثر. محاسبه خودکار مقادیر حداقل و حداکثر نیازمند جستجوی بیشتر در میان داده ها می باشد که ارزیابی و برآورد از این مقادیر ممکن است که موجب انباشتگی مقادیر نرمال شده گردد. (10 و 20 و 30) که حاصل در باز [0 و 1] بدست می آید.

$$V'(i) = \frac{(V(i) - \min(V(i)))}{(\max(v(i)))}$$

نرمال سازی انحراف معیار: نرمال سازی به روش انحراف معیار در اغلب موارد با اندازه گیری فاصله بین بازه ها بخوبی کار میکند.

برای یک مورد i مقدار مشخص با استفاده از معادله زیر تبدیل می شود: برای مجموعه $V = \{1 و 2 و 3\}$ ، $\text{mean}(V) = 2$ ، $\text{Sd}(V) = 1$ باشد و مجموعه مقادیر نرمال شده

$$V' = \{-1 و 0 و 1\}$$
 خواهد بود.

$$V'(i) = \frac{(V(i) - \text{mean}(V))}{\text{sd}(V)} \quad \text{sd} = \sqrt{S^2} \quad S^2 = \frac{\sum (x_i - x')^2}{n-1} \quad x' = \frac{\sum x_i}{n}$$

* یکنواخت سازی داده ها

یک خصیصه عددی مانند y ممکن است که بر روی مقادیر مختلفی بصورت متفاوت عمل کند این نکته را نیز باید در نظر داشت که برای بسیاری از تکنیک های داده کاوی تفاوتی بین این مقادیر با اهمیت نیستند که ممکن است با تغییر در آنها موجب کاهش ان شوند. با این وجود گاهی مواقع با یکنواخت کردن مقادیر معتبر، منجر به کاهش پیچیدگی نتایج حاصله گردد.

برای مثال یکی از روش های یکنواخت سازی داده ها، گرد کردن مقادیر با دقت معین است فرضاً برای مجموعه مقادیر معین $\{0.93, 1.01, 3.02, 2.99, 5.03, 5.01, 4.98\}$ مقادیر هموار شده برابر خواهد شد $F_{\text{smoothed}} = \{1.0, 1.0, 3.0, 3.0, 5.0, 5.0, 5.0\}$ البته الگوریتم های هموار کننده همیشه به این سادگی نیستند.

* تفاضل ها و نسبت ها

تغییرات کوچک روی مشخصه ها می توانند بهبود معنی داری را در کارایی کاوش داده ها ایجاد کند. اثرات این تبدیلات کوچک در مشخصه های ورودی/خروجی بویژه در تشخیص و کشف خطا در روش های پیش بینی کننده داده کاوی مهم هستند.

برای مثال ممکن است هدف کنترل های لازم برای فرایند تولید و رسیدن به حالت مطلوب و بهینه باشد. اما بجای نرمال سازی انحراف معیار مشخصه خروجی $s(t+1)$ ، شاید با یک حرکت نسبی از مقدار فعلی، مقدار $s(t+1) - s(t)$ در نظر گرفته شود. و یا با بکارگیری $s(t+1)/s(t)$ بعنوان خروجی فرایند داده کاوی موجب بهبود کارایی کل داده کاوی شود.

داده های از دست رفته

در بسیاری از کاربردهای دنیای واقعی کاوش داده ها، حتی با وجود مقدار داده های حجیم و فضای ذخیره سازی مناسب، ممکن است در نمونه های موجود، مقادیری از داده ها از دست

رفته (گمشده) باشند. در بعضی روش های داده کاوی، مقادیر از دست رفته داده ها و فرایند داده ای مناسب برای رسیدن به نتیجه نهایی را می پذیرند (البته این نوع داده کاوی ها بیشتر برای تعداد مقادیر کوچک از دست رفته مصداق دارد)

اما مشکل از آنجا آغاز میشود که برای مجموعه داده های بزرگ نمی توان از مقادیر از دست رفته چشم پوشی کرد. یک راه حل برای جایگزینی خودکار مقادیر از دست رفته با مقادیر ثابت عبارت است از:

1- جایگزینی تمام مقادیر از دست رفته با یک تک مقدار ثابت سراسری

2- جایگزینی یک مقدار از دست رفته با متوسط مشخصه آن.

3- جایگزینی یک مقدار از دست رفته با متوسط مشخصه آن برای یک گروه

مشخص.

نکته ای که باید به آن توجه کرد این است که در همه روش های بالا اشکال در این است که مقدار جایگزین شده مقدار درست و واقعی آن نیست که در این حالت به مقادیر نزدیک به مقدار از دست داده شده، یک کلاس همگن یا کلاس مصنوعی گفته می شود.

تحلیل داده های نامنطبق

در مجموعه داده های بزرگ، به نمونه هایی که از رفتار کلی مدل داده ای تبعیت نمی کنند و بطور کلی متفاوت یا ناهماهنگ با مجموعه باقیاننده داده ها هستند، داده های نامنطبق نامیده می شوند.

داده های نامنطبق می توانند توسط خطای اندازه گیری ایجاد شوند یا نتیجه نوع داده ای درونی باشند.

برای مثال اگر سن فردی در پایگاه داده 1- باشد، مقدار فوق قطعاً غلط با یک مقدار پیش فرض فیلد "سن ثبت نشده" می تواند در برنامه مشخص گردد. و نمونه های دیگری که ممکن است بر اثر خطای باشند...

و اما راه حل بسیار از الگوریتم های کاوش داده ها سعی در کاهش وبعضا حذف داده های نامنتطبق در مراحل پیش پردازش دارند.

نکته در اینجا تحلیلگران داده کاوی باید در الگوریتم های پیش پردازنده پیش از حذف خودکار داده های نامنتطبق، بر جنبه های آشکارسازی داده های نامنتطبق تاکید کنند.

روش های آماری یکی از ساده ترین راه ها برای تشخیص وکشف داده های نامنتطبق روش های آماری است.

بافرض اینکه توزیع مقادیر داده ای انجام شوند، یافتن پارامترهای آماری اساسی مثل مقدار متوسط و انحراف معیار لازم می شوند. براساس همین مقادیر و تعداد داده های نامنتطبق، ایجاد مقدار استانه به عنوان یک تابع آماری بوجود می آید.

مثال :مجموعه داده های مشخصه ای مثل سن، با مقدار 20:

Age={37و45و55و31و439و28و9و22و41و52و156و39و23و56و3}

میانگین =39.9

انحراف معیار = 45.65

انحراف معیار *2 ± میانگین =مقدار آستانه یا فاصله

همه داده هایی که خارج از بازه [131.2و-54.1]قرار دارند داده های نامنتطبق هستند. که با شناخت ویژگی

های مشخصه (سن همیشه بزرگتر از صفر) ممکن است که باعث کاهش بیشتر بازه گردد[131و0]در مثال

ما مقدار داده های نامنتطبق در پایه شاخص معین، 154و139و-67 می باشند.

تشخیص داده های نامنتطبق بر مبنای فاصله روشی دیگر است که محدودیت های ابعاد داده ها در روش آماری را ندارد.

روش ها و تکنیک های بر مبنای انحراف سومین کلاس از روش های تشخیص داده های نامنطبق، دراصل شبیه سازی روشی است که می توان نمونه های غیر طبیعی را از یک مجموعه نمونه های مشابه دیگر تشخیص داد.

این روش ها ویژگی های اساسی مجموعه نمونه ها را تعیین و نمونه های متفاوت از این ویژگی ها منحرف می شوند و به تبع آن داده های منطبق آشکار میشوند. البته باید توجه داشت که تعیین ویژگی های مجموعه نمونه ها بسیار پیچیده و گاهی ترکیبی از انتخاب طرح های کلی تر می باشند.

فصل سوم

وب کاوی

۱-۳ وب کاوی

استفاده از وب داده های وب یکی از گام های کلیدی در کشف دانش در پایگاه داده، ایجاد یک مجموعه داده مناسب جهت انجام داده کاوی می باشد. در وب کاوی این داده می تواند از سمت سرور، مشتری، پروکسی سرور یا از یک پایگاه داده سازمان جمع آوری شود.

هر کدام از این داده ها نه تنها از نظر منابع داده متفاوت می باشند بلکه از نظر انواع داده های موجود و محدوده مکانی که آن داده از آنجا جمع آوری می شود و متد پیاده سازی انواع داده ای که در وب کاوی استفاده می شود شامل:

محتوا: داده واقعی در صفحات وب، داده ای که صفحه وب برای نمایش آن به کاربران طراحی شده است. که معمولاً از متن و گرافیک تشکیل شده ولی به آن محدود نمی شود.

ساختار: داده ای که سازمان دهی محتوا را مشخص می سازد. اطلاعات ساختار درون صفحات شامل ترتیب انواع تگ های XML یا HTML در یک صفحه داده شده می باشد و می تواند به صورت یک ساختار درختی نمایش داده شود که تگ ریشه درخت می باشد. اصلی ترین نوع از اطلاعات ساختاری بین صفحات، هایپرلینک^۱ است که یک صفحه را به دیگری مرتبط می کند.

استفاده: داده ای که الگوی استفاده از صفحات وب را مشخص می سازد، مثل آدرس های IP، رجوع به صفحات و تاریخ و زمان دسترسی پروفایل

کاربر: داده ای که اطلاعات آماری درباره کاربران وب سایت فراهم می سازد که شامل داده ثبت نام و اطلاعات پروفایل مشتری می باشد.

منابع داده داده های استفاده که از منابع مختلفی جمع آوری می شود، الگوهای راهبری از بخش های مختلفی از کل ترافیک وب را نمایش می دهد. جمع آوری در سطح سرور لاگ های وب سرور یک منبع مهم برای اجرای وب کاوی استفاده از وب محسوب می شود زیرا به طور صریح رفتار مرورگری تمام مشاهده کنندگان سایت را ثبت می کند. داده ای که در لاگ سرور ثبت می شود، دسترسی به یک وب سایت که از سوی تمام کاربران صورت می گیرد را منعکس می کند. این فایل های لاگ به فرمت های گوناگونی چون Common log یا Extended log ذخیر می شوند.

^۱ hyperlink

جمع آوری در سطح مشتری جمع آوری داده در سطح مشتری می تواند با بکارگیری یک عامل از راه دور (مثل اپلت های جاوا یا جاوا اسکریپت) یا با تغییر کد مرجع یک مرورگر موجود (مثل Mozilla یا Mosaic) پیاده سازی شود. پیاده سازی این نوع روش جمع آوری داده در سطح مشتری به همکاری کاربر در هر دو مورد ذکر شده نیاز دارد. جمع آوری در سطح پروکسیک پروکسی وب به عنوان یک سطح میانی از ذخیره سازی بین مرورگر سمت مشتری و وب سرور محسوب می شود تا زمان بارگذاری صفحه وبی که توسط کاربر تجربه شده را کاهش دهد همانطور که بار ترافیکی در سمت مشتری و سرور را کاهش می دهد. داده های لاگ مربوط به وب معمولاً حجیم و گسترده هستند و به منظور کشف الگو، این داده ها باید در یک دید یکپارچه، سازگار و جامع جمع آوری شوند. در بیشتر کاربردهای داده کاوی پیش پردازش داده با حذف و فیلتر کردن داده های افزونه و بی ربط و حذف نویز و تبدیل و رفع هر ناسازگاری سروکار دارد. پیش پردازش داده نقش اساسی در کاربردهای کشف دانش در داده استفاده از وب دارا هستند و مهمترین مساله در بیشتر روش های کشف الگو، مشکل آن ها در اداره داده های استفاده از وب در مقیاس بزرگ است. به همین خاطر اکثر فرایندهای KDWUD به طور غیر بر خط انجام می شوند.

تحلیل داده استفاده بدون روش پیش پردازش مناسب نتایج ضعیف و یا حتی خرابی را بدنبال خواهد داشت. بنابراین متودولوژی برای پیش پردازش باید به کار گرفته شود تا هر مجموعه ای از فایل های لاگ وب سرور را به مجموعه ساختاریافته ای از جداول در مدل پایگاه داده رابطه ای تبدیل کند.

فایل های لاگ از روی سایت های مختلف یک سازمان با هم ادغام می شوند تا رفتار کاربرانی که از طریقی ملموس راهبری داشته اند را نمایش دهد. بنابراین این فایل ها باید با حذف درخواست هایی که مورد نیاز نیستند، پاک می شوند مانند درخواست های ضمنی برای آبجکت های تعبیه شده در صفحات وب و یا درخواست هایی که بوسیله مشتری های غیر انسانی وب سایت ها ایجاد می شود. درخواست های باقیمانده با کاربر، نشست های کاربر و مشاهدات و مشاهده صفحات، گروه بندی می شود. و در نهایت مجموعه های پاک و تبدیل شده از درخواست های کاربران در یک مدل پایگاه داده رابطه ای ذخیره می شود.

از فیلتر هایی برای فیلتر کردن داده های بدون استفاده، بی ربط و ناخواسته استفاده می شود تحلیلگر می تواند فایل های لاگ را از وب سرورهای متفاوت جمع آوری کند و تصمیم گیری کند که کدامیک از ورودی ها مطلوب هستند. در واقع هدف این است که اندازه بزرگ داده های استفاده از وب موجود به طور قابل توجهی کاهش یابد و در عین حال کیفیت آن با سازمان دهی آن و فراهم سازی متغیر های یکپارچه اضافی برای تحلیل داده کاوی افزایش یابد.

فرمول بندی مسائل فرض کنید مجموعه $R = \{r_1 + r_2 + \dots + r_m\}$ تمام منابع وباز یک وب سایت باشد. اگر U مجموعه تمام کاربرانی که به سایت دسترسی داشتند، باشد؛ عنصر لاگ بصورت li تعریف می شود که $ui \in U; ri \in R$ است و t زمان دسترسی را نمایش می دهد، s وضعیت درخواست و ref آصفحه مورد مراجعه را نمایش می دهد ref_i . در برخی از فرمت های لاگ های وب مثل CLF در حالی که صفحه مورد مراجعه ثبت نشده، اختیاری است s . یک کد سه رقمی است که موفقیت یا شکست درخواست مورد نظر را نشان می دهد. همچنین در موارد دیگر دلیل شکست را نیز بیان می کند. یک وضعیت با مقدار $s=200$ نشان می دهد که درخواست موفق است در حالی که وضعیت با مقدار $s=404$ نشان دهنده این است که فایل مورد درخواست در محل مورد نظر یافت نشده است. $li = \{li_1, li_2, \dots, li_m\}$ به ترتیب صعودی ذخیره می شوند که یک لاگ وب سرور را تشکیل می دهند. در صورت داشتن N وب سرور، مجموعه فایل های لاگ $Log = \{L_1, L_2, \dots, L_N\}$ است. با بکارگیری این علائم مسئله پیش پردازش به صورت زیر فرمول بندی می شود.

" با دریافت یک مجموعه از فایل های لاگ مربوط به لاگ های وب سایت های مختلف، کاربر، نشست های کاربر، مشاهده و مشاهدات صفحات کاربران وب سایت در یک بازه زمانی مشخص Δt استخراج می شود."

پیش پردازش داده

فرایند پیش پردازش گام های زیر را در بر می گیرد: ادغام فایل های لاگ از وب سرورهای گوناگون پاک کردن داده شناسایی کاربران، نشست ها و مشاهده هافرمت بندی داده و خلاصه سازی آن ادغام در ابتدای پیش پردازش داده، درخواست از تمام فایل های لاگ در Log در یک فایل لاگ الحاقی همراه با نام وب سرور جهت تشخیص بین درخواست های ایجاد شده مربوط به وب سرورهای مختلف و همچنین توجه به همگام سازی کلاک های وب سرورهای مختلف که از لحاظ زمانی متفاوت اند.

به خاطر دلایل محرمانگی، فایل لاگ نتیجه f را بی نام کرده بطوریکه وقتی که فایل های لاگ به اشتراک گذاشته می شود یا نتایج منتشر می شوند، نام میزبان یا آدرس های IP، از بین می روند. بنابراین نام اصلی میزبان با یک شناسنده ای که اطلاعاتی درباره محدوده دامنه کد کشور یا نوع سازمان مثل edu، .org، .com) نگهداری می کند، جایگزین می شود.

مسئله ادغام به صورت زیر فرمول بندی می شود: با دریافت یک مجموعه فایل های لاگ $Log = \{L_1, L_2, \dots, L_n\}$ این فایل های لاگ در یک فایل لاگ مجزا و منفرد ادغام می شود (فایل لاگ

الحاقی (فرض کنید Li ، آیین فایل لاگ می باشد $Li.c$. را به عنوان اشاره گر بر روی درخواست های Li در نظر بگیرید و $Li.1$ عنصر لاگ جاری از Li است که با $Li.c$ نشان داده می شود و $Li.1.time$ ، زمان t مربوط به عنصر لاگ جاری از Li می باشد و همچنین $S=(w_1, w_2, \dots, w_n)$ آرایه ای از اسامی وب سرورها می باشد به طوری که $S[i]$ نام وب سرور مربوط به لاگ $Li.1$ می باشد .

مراحل :

- 1- مقداردهی اولیه اشاره گر فایل لاگ الحاقی ξ
- 2- اسکن عناصر لاگ از هر فایل لاگ Li در Log و افزودن آن به ξ
- 3- مرتب سازی عناصر ξ به طور صعودی بر اساس زمان دسترسی آن ها برگرداندن مقدار ξ پاک کردن داده .

گام دوم در پیش پردازش داده حذف درخواست های بدون استفاده از فایل های لاگ می باشد . بطوریکه اگر تمام ورودی های لاگ معتبر نباشند، باید ورودی های بی ربط را حذف کنیم . معمولاً این فرایند تمام درخواستهایی که منابع غیر قابل تحلیل مثل تصاویر، فایل های چندرسانه ای و فایل های مربوط به سبک صفحات را در بر می گیرند، را حذف می کند .

برای مثال درخواستهای مربوط به محتوای صفحات گرافیکی تصاویر ($*.jpg$ & $*.gif$) و همچنین درخواستهای مربوط به هر فایل دیگر در یک صفحه وب یا حتی نشست های راهبری که توسط ربات ها و اسپایدر های وب انجام می شود .

با فیلتر کردن داده های بی استفاده، می توانیم سایز فایل لاگ را کاهش داده تا از فضای ذخیره سازی کوچکتری استفاده کرده و نیز کارهای بعدی را آسان تر کنیم . برای نمونه، با فیلتر کردن درخواست های تصاویر، سایز فایل های لاگ وب سرور نسبت به سایز اولیه اش تا 50 درصد کاهش می یابد . بنابراین پاک کردن داده حذف ورودی های بی ربطی چون موارد زیر می باشد:

1- درخواستهایی که توسط برنامه های خودکار انجام می شود مثل $\text{Web Robot, Spiders}$ و Crawler ها . این برنامه ها ترافیکی بر روی وب سایت ایجاد می کنند که می توانند بر روی آمار سایت تاثیر بگذارند و همچنین در بررسی هایی که توسط KDWUD انجام می شود مطلوب نیستند .

2- درخواستهای مربوط به فایل های تصویری که به صفحات مشخصی اختصاص داده می شود .

3- درخواست یک کاربر برای مشاهده یک صفحه خاص معمولاً در چندین در چندین عنصر از لاگ منعکس می شود زیرا هر صفحه گرافیک هایی را شامل می شود که فقط آنهایی برای ما مهم هستند که کاربر صریحاً آنها را درخواست کرده که معمولاً فایل های متنی هستند .

عناصر با کدهای وضعیت HTTP نا موفق . کدهای وضعیت HTTP برای نشان دادن موفقیت یا شکست یک درخواست بکار می روند که در اینجا ما فقط عناصر با کد بین 200 تا 299 که با موفقیت انجام شده اند در نظر می گیریم .

عناصری که متدی به غیر از GET و POST دارند شناسایی در این گام درخواستهای غیر ساختیافته یک فایل لاگ به صورت کاربر (user)، نشست کاربر (user session)، مشاهدات و ملاقات صفحات (page view, visit) گروه بندی می شود . در پایان این گام فایل لاگ به صورت یک مجموعه از تراکش ها خواهد بود (نشست کاربر یا مشاهدات) کاربرد بیشتر موارد فایل لاگ فقط آدرس های کامپیوتر نام یا IP و عامل کاربر را فراهم می سازد به عنوان مثال فایل های لاگ . ECLF برای وب سایتی که نیازمند ثبت کاربر هستند، فایل لاگ همچنین User login را شامل می شود (به عنوان سومین رکورد در یک عنصر لاگ) که برای شناسایی کاربر استفاده می شود . وقتی که user login موجود نباشد هر IP به عنوان کاربر در نظر گرفته می شود . با این حال این واقعیت وجود دارد که یک آدرس IP توسط چندین کاربر استفاده می شود و این برای KDWUD جهت شناسایی کاربر کافی نیست . به هر حال هنوز هم مکانیزمی برای تشخیص و تمایز بین کاربران برای تحلیل رفتار دسترسی کاربر مورد نیاز است .

نشست کاربر شناسایی نشست کاربر از فایل لاگ بدلیل پروکسی سرورها، آدرس های پویا و مواردی که چندین کاربر از طریق یک کامپیوتر دسترسی پیدا می کنند (در کتابخانه، کافی نت و...) یا یک کاربر از چندین مرورگر یا کامپیوتر استفاده می کند، امکان پذیر نمی باشد . یک نشست کاربر به صورت ترتیبی از درخواست ها که بوسیله یک کاربر منفرد در یک دوره زمانی مشخص تعریف می شود . یک کاربر می تواند یک (یا چند) نشست در طول یک دوره زمانی داشته باشد . شناسایی نشست عبارت است از فرایند قطعه بندی لاگ دسترسی هر کاربر به نشست های دسترسی مجزا . دو روش بر اساس زمان وجود دارد که شامل روش مبتنی بر طول نشست (Session-duration) و روش مبتنی بر page-stay-time همچنین می توانیم از یک آستانه زمانی timeout استفاده می کنیم .

۲-۲ روش‌های وب‌کاوی

داده‌های موجود در وب، انواع مختلفی دارند. این انواع داده‌ها را می‌توان به سه گروه کلی تقسیم کرد:

- 1- محتوای وب Web Content ، مانند متن، تصاویر، جداول و
- 2- کاربرد وب Web Usage ، مانند فایل‌های LOG بدست آمده در سرورها و
- 3- ساختار وب Web Structure ، مانند لینک‌ها و برچسب‌ها.

حال با توجه به تقسیم‌بندی بالا، وب‌کاوی به سه دسته کلی ذیل تقسیم می‌گردد:

- 1- وب‌کاوی محتوا Web Content Mining
- 2- وب‌کاوی کاربرد Web Usage Mining
- 3- وب‌کاوی ساختار Web Structure Mining

در سه بخش پیش‌رو به بررسی مختصر هر یک از این دسته‌ها خواهیم پرداخت .

Web Content Mining وب‌کاوی محتوا، فرآیند استخراج اطلاعات کاربردی از محتوای وب را گویند. محتوای وب، داده‌های ساخت‌یافته (Structured) مانند جدول، نیمه ساخت‌یافته مانند برچسب‌های HTML و غیرساخت‌یافته (Unstructured) مانند متن معمولی را شامل می‌شود. این حیطه از روش‌های دیگری مانند داده‌کاوی متن (Text Mining) ، پردازش زبان طبیعی (NLP) و بازیابی اطلاعات (Information Retrieval) در دل خود استفاده می‌نماید. در وب‌کاوی محتوا تکنیک‌های داده‌کاوی معمولی مثل Classification ، Clustering و Association را نیز می‌توان به کار گرفت .

Web Usage Mining استخراج الگوهای معنادار از داده تولید شده در سرورها، پست الکترونیکی و ... را وب‌کاوی کاربرد گویند. با کمک این روش می‌توان رفتار کاربران مختلف را پیش‌بینی کرده و بر اساس آن، صفحات وب را برای آن کاربر خاص، شخصی‌سازی نمود. با استفاده از وب‌کاوی کاربرد می‌توان به شخصی‌سازی صفحات وب پرداخت (Personalization) ، سایت‌ها را اصلاح نمود (Site Modification) و یا کرائی سایت‌ها افزایش داد

نکته دیگر در باره فرمت فایل‌های LOG است. این فایل‌ها فرمت استاندارد خاصی ندارند و در هر کاربرد و هر سرور بر حسب نیاز خود، فرمت و قالب مخصوصی را تعریف کرده و از آن قالب خاص استفاده می‌کنند.

Web Structure Mining وب‌کاوی ساختار به استخراج اطلاعات از ساختار صفحات وب می‌پردازد. این روش خود به دو دسته‌ی Intra-page و Inter-page تقسیم می‌شود.

در مورد روش Intra-page ، ساختار درونی یک صفحه وب، در نظر گرفته می‌شود، در صورتی که در روش Inter-page ساختار صفحات وب در ارتباط با یکدیگر مورد توجه است. به روش Inter-page ، تحلیل ابرلینک‌ها (Hyperlink Analysis) نیز گفته می‌شود. در این روش معمولاً صفحات وب و لینک‌های میان آنها به شکل گراف جهت‌دار نمایش داده می‌شود. هر گره از این گراف، نشان‌دهنده یک صفحه وب بوده و هر یال جهت‌دار از گره "الف" به گره "ب" وجود لینکی از صفحه "الف" به صفحه "ب" را نمایش می‌دهد.

فصل چهارم

مدیریت ارتباط با

مشتری

۴-۱ تاریخچه:

شاید بتوان تاریخچه ظهور مباحث مرتبط به CRM را در سه دوره زیر خلاصه نمود:

الف) دوره انقلاب صنعتی (تولید دستی تا تولید انبوه) ابتکارات فورد در بکارگیری روش تولید انبوه به جای روش تولید دستی، یکی از مهمترین شاخص‌های این دوره می باشد. هر چند تغییر شیوه تولید باعث شد که محدوده انتخاب مشتریان از نظر مشخصه‌های محصول کاهش یابد (نسبت به تولیدات صنایع دستی) اما محصولات تولید شده به روش جدید از قیمت تمام شده پایین تری برخوردار شدند. به عبارتی دیگر در انتخاب روش تولید انبوه از سوی فورد، افزایش کارایی و صرفه اقتصادی مهمترین اهداف پیش بینی شده بودند.

ب) دوره انقلاب کیفیت (تولید انبوه تا بهبود مستمر) این دوره هم‌زمان با ابتکار شرکت‌های ژاپنی مبنی بر بهبود مستمر فزایندها آغاز شد. این امر به نوبه خود به تولید کم هزینه تر و با کیفیت تر محصولات منجر شد. با مطرح شدن روشهایی نوین مدیریت کیفیت مانند TQM این دوره به اوج خود رسید. اما با افزایش تعداد شرکت‌های حاضر در عرصه رقابتی و گسترش فرهنگ حفظ و بهبود کیفیت محصول (از طریق ابزارهای مختلف کیفیتی) دیگر این مزیت رقابتی برای شرکتها پیشرو و کارساز نبوده و لزوم یافتن راه‌های جدیدی برای حفظ مزیت رقابتی احساس می شد.

ج) دوره انقلاب مشتری (بهبود مستمر تا سفارشی‌سازی انبوه) در این دوره با توجه به افزایش توقع مشتریان، تولید کنندگان ملزم شدند محصولات خود را با هزینه کم، کیفیت بالا و تنوع زیاد تولید کنند. به معنای دیگر تولید کنندگان مجبور بودند توجه خود را از تولید صرف به یافتن راه‌هایی برای رضایت مشتریان سابق خود معطوف نمایند.

۲-۴ مشتری کیست

کسی که محصولات یا خدمات سازمان را خریده و یا از آن استفاده میکند . به عنوان مشتری تعریف میگردد . در مجموعه میتوان مشتریان را به دو دسته تقسیم کرد . مشتریان خارجی و مشتریان داخلی مشتریان خارجی در بیرون از سازمان بوده و محصولات و خدمات میخرند . تمام کارمندانی که به نوعی در فرآیند های تهیه و توضیح کالا یا خدمات نقش دارند باید اثر کار و وظیفه خود را در بالابردن سطح رضایت مشتریان خارجی به خوبی درک نمایند .

در برابر مشتریان خارجی هر سازمانی مشتریان داخلی دارند که به اندازه مشتریان خارجی مهم اند . در تمام مراحل از قسمت های مهندسی گرفته تا تولید و دیگر فرآیندها همواره یه مشتری داخلی وجود دارد که محصول و خدمات را دریافت میکنند و در عوض محصول خدمتی را ارایه مینمایند

۳-۴ 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.

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.

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.

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

۴-۳ انواع CRM:

1-3-4 CRM عملیاتی (Operational):

در این روش کلیه مراحل ارتباط با مشتری، از مرحله بازاریابی و فروش تا خدمات پس از فروش و اخذ بازخورد از مشتری، به یک فرد سپرده می شود، البته به نحوی که فروشندگان و مهندسان ارائه خدمات بتوانند سابقه هر یک از مشتریان را بدون مراجعه به این فرد در دسترس داشته باشند. از ابزار و روشهای CRM عملیاتی میتوان به SFA یا قدرت فروش مکانیزه اشاره نمود که کلیه عملیات به مدیریت تماس بورس و مدیریت اداره فروش را برعهده دار CSS ابزار دیگر CRM عملیاتی است که در آن به جای ارتباط تلفنی با مشتری، از ابزارهای دیگری مانند ارتباط رو در رو، اینترنت، فاکس و کیوسکهای مخصوص پاسخگویی به مشتریان استفاده می شود.

2-3-4 CRM تحلیلی (Analytical):

در CRM تحلیلی ابزارها و روشهایی به کار می رود که اطلاعات به دست آمده از CRM عملیاتی را تجزیه و تحلیل نموده و نتایج آن را برای مدیریت عملکرد تجاری آماده می کند. این سیستم مهمترین نوع از CRM می باشد. به این صورت که شامل دادههایی است که برنامهها

جهت برقراری ارتباط با مشتری به آن نیاز دارند. به عبارت دیگر این داده‌های خام در اختیار برنامه‌های CRM قرار می‌گیرند و پس از کار بر روی این داده‌ها، نتیجه مناسب در اختیار شرکت و مشتری قرار داده می‌شود. اما اگر بخواهیم یک تعریف کامل ارائه نماییم: بدست آوردن، ذخیره، پردازش، تفسیر و ارائه گزارش به استفاده‌کنندگان داده‌ها مشتری می‌باشد. شرکتهای زیادی هستند که این داده‌ها را جمع‌آوری کرده و پس از استفاده از الگوریتمهایی مختلف سعی در تحلیل و تفسیر این داده‌ها می‌نمایند. در واقع، CRM عملیاتی و تحلیلی در یک تعامل دو طرفه هستند.

4-3-3 CRM تعاملی (Collaborative):

در این نوع ارتباط، مشتری برای برقراری ارتباط با سازمان از سهل‌ترین روش ممکن مانند تلفن، تلفن همراه، فاکس، اینترنت و سایر روش‌های مورد نظر خود استفاده می‌نمایند. این نرم‌افزارها را PRM می‌نامند. CRM تعاملی به دلیل امکان انتخاب روش از سوی مشتری و اینکه اکثر فرآیندها (از جمع‌آوری داده‌ها تا پردازش و ارجاع مشتری) در حداقل زمان ممکن به مسوول مربوطه صورت می‌گیرد باعث مراجعه مجدد مشتری و ادامه ارتباط با شرکت می‌شود

4-4 نقش مدیران ارشد در ارتباط با مشتریان

یک شرکت هوپیمایی را در نظر بگیرید. اگر این شرکت به لحاظ ارتباط با مشتریان ضعیف باشد، به آسانی می‌توان آن را زیر سوال برد و از آن انتقاد کرد.

جان تگ، مدیر ارشد عملیاتی شرکت هوپیمایی یونایتد ایرلاینز قسمت اعظم کار خود را روی ارتقای سیستم‌های اطلاع‌رسانی به مشتری، تکنولوژی‌های مشتری‌مدار و قابلیت‌های خدمات به مشتری متمرکز کرده است. وی در مصاحبه‌ای بیان داشت «بزرگ‌ترین چالش ما آگاه کردن سازمان به میزان فعالیت لازم برای ارتباط بهتر با مشتری است. بدون احساس ضرورت، حقیقتاً تغییر دادن فرهنگ مان بسیار سخت است. حال که یونایتد ایرلاینز با شرکت هوپیمایی Continental ادغام شده است امید آن می‌رود که قدرت ارتباطی آن بیشتر شود. دلیل دوم که چرا مدیران شرکتهای بزرگ ارتباط محدودی بامشتریان دارند، این است که آنها قدرت نمادین روش‌های مدیریتی خود را دست کم یا نادیده می‌گیرند. نمونه و الگو بودن جزو ماهیت ثانویه مدیران کارفرما است و مدیران ارشد شرکتهای بزرگ باید نسبت به این جنبه از مدیریت آگاهی بیشتری داشته باشند. مدیران ارشد نمی‌توانند وظیفه فرهنگ‌سازی در سازمان را به فرد

دیگری واگذار کنند. آنها موثرترین و مهم‌ترین شخص در شرکت‌ها هستند که قادر به انجام این کار هستند. شیوه‌های مدیریتی مدیران ارشد روی عملکرد تیم‌های کاری، کارمندان و مدیران دیگر تاثیر می‌گذارد. زمانی که ارتباط مدیران ارشد با مشتریان قطع شود برای کل سازمان نیز چنین اتفاقی رخ می‌دهد. سومین دلیلی که ارتباط با مشتری به مساله مبهمی تبدیل می‌شود این است که به عنوان یک مساله مهم مدیریتی در نظر گرفته نمی‌شود. امروزه دغدغه و نگرانی مدیران مدرن مسائلی چون استراتژی، ادغام و تملک، قیمت سهام شرکت و قدرت ترازنامه است. ارتباط با مشتری کار سطح پایینی تلقی می‌شود و آنقدر پیش پا افتاده در نظر گرفته می‌شود که نمی‌توان آن را جزء اصلی فرمول موفقیت سازمان در نظر گرفت. چهارمین و آخرین دلیلی که شرکت‌های بسیاری وجود دارند که ارتباط مناسبی با مشتریان ندارند این است که اکثر مدیران از تمام جنبه‌ها به مشتریان فکر نمی‌کنند. از دید یک مدیر ارشد در شرکتی بزرگ، مشتریان معمولاً اطلاعات، نمودار یا اعداد و ارقام هستند. آنها به عنوان انسان‌های واقعی دیده نمی‌شوند. نقطه کور مدیریت کوتاهی در ایجاد سازمان‌های مشتری‌مدار و عدم درک نقشی که مدیریت ایفا می‌کند اشتباه بزرگ و هزینه‌برداری است. عدم ارتباط حقیقی و پایدار منجر به نابودی عواملی چون نوآوری، وفاداری مشتری و رشد سودآور در شرکت‌ها شده است. مدیران اندکی دریافته‌اند که چگونه سازمان‌های مشتری‌مدار را برپا کنند و فقدان چنین مدیریتی میلیاردها دلار سرمایه را در اقتصادهای جهانی و آمریکا از بین برده است. بارزترین مثال، صنعت خودروسازی ایالات متحده است که به دلیل عدم ارتباط مناسب مدیران با مشتریان فلج شد. عدم برقراری ارتباط با مشتری برای بسیاری از صنایع دیگر نیز تبدیل به مساله بزرگی شده است. شرکت‌های کرایه اتومبیل در حال تلاش برای متعادل کردن کاهش هزینه‌ها و خدمات به مشتری هستند. این وضعیت در مورد بیمارستان‌ها، شرکت‌های هواپیمایی و بانک‌ها نیز همین‌طور است. تعداد اندکی در این صنایع دارای مدیران مشتری‌مدار هستند و آنهایی که ارتباط مناسبی با مشتریان دارند (برای مثال شرکت هواپیمایی Southwest به داشتن مشتریان وفادار، تقویت نوآوری و لذت بردن از کسب و کار پایدار مشهور هستند). نادیده گرفتن اهمیت ارتباط با مشتری قابل توجه دیگری نیز داشته است. سال 2008 در زمان رکود اقتصادی شرکت‌های آمریکایی بیش از 130 میلیارد دلار برای آموزش و رشد کارمندان هزینه کردند که از این میان، حدود 50 میلیارد دلار صرف ارتقا و پیشرفت مدیران شد. این میزان مبلغی است که شرکت‌ها

هزینه کردند تا از داشتن مدیران مستعد و مجرب برای آینده‌ای رو به رشد اطمینان حاصل کنند. یکی از رایج‌ترین و موثرترین راه‌ها برای پرورش مدیران آینده‌سنجیدن مهارت‌های مدیریتی آنها و ارائه بازخورد به مدیران است؛ یعنی اینکه به آنها نشان دهد تا چه میزان این مهارت‌ها را دارا هستند. به منظور انجام این کار، شرکت‌ها باید دقیقاً تصمیم بگیرند چه مهارت‌هایی مهم‌تر هستند و ارزش ارزیابی دارند. این لیست از مهارت‌ها اغلب به عنوان «برنامه موفقیت» تلقی می‌شوند. راه‌حل چیست؟ در اینجا به ذکر چند راهکار برای ایجاد ارتباط مناسب با مشتریان می‌پردازیم. جای هیچ تعجبی نیست که این راهکارها همان مشخصه‌های مدیران موفق است: به ایجاد تنظیمات سازمانی جدید پردازید. در شرکت‌های کوچک، تشکیلات سازمانی ساده است و برای مدیران ارتباط نزدیک با مشتری آسان است. در شرکت‌های بزرگ‌تر، تشکیلات سازمانی دارای پیچیدگی است. نظام سلسله‌مراتب و کاغذ بازی مانع رسیدن صدای مشتری به مدیران می‌شود. مدیران ارشد به دلیل نداشتن دید روشن قادر به تشخیص یا پاسخ سریع به مشتریان نیستند. در شرکت‌های متوسط، یک راه موثر برای از بین بردن موانع سازمانی، حذف دیوارها به طور فیزیکی بین سطوح است. شرکت Zappos با بیش از دوهزار کارمند دست به چنین اقدامی زد: در دفتر مرکزی پارتیشن‌هایی که میزکار کارکنان را از قسمت‌های دیگر مجزا می‌کند به حدی کوتاه است که افراد می‌توانند بایستند و به راحتی با آنها صحبت کنند. شرکت Communispace نیز یکی از سازمان‌هایی است که به ارتباط با مشتری بسیار اهمیت می‌دهد. دفتر مرکزی این شرکت در ایالت ماساچوست آمریکا دارای طراحی فضای باز است. خانم دیان حسان، مدیر ارشد این شرکت طی مصاحبه‌ای گفت: نمی‌توانم شکل دیگری از محیط کار را تصور کنم. این مدل با فرهنگ، ارزش‌ها و تعهد ما به شفافیت و قابل دسترس بودن مطابقت دارد. اقدامات مشتری‌مداری را ارزیابی کنید و برای آن پاداش قائل شوید. سیستم‌های ارزیابی و پاداش، چه رسمی و چه غیررسمی کمک می‌کنند تا تعیین کنید افراد به چه چیزی توجه می‌کنند و چگونه رفتار می‌کنند. مدیران ارشد باید در دادن پاداش به افرادی که برای ارتباط با مشتری تلاش بسیار می‌کنند پیشقدم باشند. یک گفته مشهور می‌گوید «هرآنچه که ارزیابی شود، مدیریت می‌شود» و این دستورالعمل در هر جنبه‌ای از ارتباط با مشتری کاربرد دارد. مساله اصلی این است که چه چیز باید مورد سنجش و ارزیابی قرار بگیرد و این ارزیابی چگونه باید انجام شود. در صورت امکان، معیارهای سنجش باید فراتر از داده‌های کمی و آماری سنتی

حاصل از تحقیقات بازار باشد. مدیران ارشد باید اصرار داشته باشند که صدای مشتری شنیده شود و کوچک‌ترین جزئیات آن شرح داده شود. داده‌های کمی و آمار و ارقام مفید هستند اما معیارهای اندازه‌گیری کیفی و به صورت حکایت و داستان دارای عمق بیشتری است. صحبت و گفت‌وگو بیشتر از اعداد و ارقام به ما اطلاعات می‌دهند. از تکنولوژی آنلاین استفاده کنید. با توجه به اندازه، تنوع و ماهیت جوامع مشتریان، ارتباط با آنها به شیوه‌ای اصولی و متفکرانه بدون استفاده از تکنولوژی آنلاین تقریباً غیرممکن است. اینترنت این امکان را می‌دهد که سازمان‌ها، اطلاعات بسیار زیاد و همین‌طور ارزشمندی را جمع‌آوری کنند. مدیران ارشد شرکت‌ها به جای عقب کشیدن و قطع ارتباط با میلیون‌ها مشتری باید از فناوری اینترنت استفاده کنند و اصرار داشته باشند که سازمان‌ها برای بهترشنیده شدن صدای مشتری از این تکنولوژی بهره ببرند. همچنین تکنولوژی آنلاین برای مدیران پرمشغله این امکان را فراهم می‌کند تا بر موانع زمانی غلبه کنند. ضرورت فرهنگی سازمان‌هایی که بامشتریان خود ارتباط دارند، جو سازمانی به‌خصوصی دارند. این موضوع اصلاً ربطی به اندازه یا نوع فعالیت شرکت ندارد، چنین محیط‌های کاری دارای ویژگی‌های مهم مشترکی هستند: آنها خیلی ساختارگرا نیستند؛ دارای استاندارد بالا هستند؛ برای مسوولیت فردی اهمیت بسیاری قائلند؛ روحیه همکاری بالایی دارند؛ دارای «سیستم باز» هستند یعنی با محیط اطراف خود تعامل دارند و در آخر تعهد بسیار بالایی به شرکت، اهداف آن و کارمندان دارند. مدیران چگونه باید این جو سازمانی را ایجاد کنند؟ برای شکل‌گیری فرهنگ ارتباط با مشتریچه اقدامات مدیریتی مورد نیاز است؟ فرهنگ‌سازی صحیح وابسته به برخی اقدامات مدیریتی است که با محیط سازمانی دارای عملکرد بالا و موفقیت کسب و کار همراه شده است. به چند مورد از این اقدامات اشاره‌ای خواهیم کرد:

1- اهداف چالش‌انگیز را برای ارتباط با مشتریان قرار دهید. این به معنی وادار کردن سازمان‌ها به اتخاذ اهداف بزرگ و به چالش کشیدن سازمان برای بهبود ارتباط با مشتری است

2- تعهد شخصی خود را برای دستیابی به اهداف ایجاد ارتباط با مشتری نشان دهید. این بدان معنا است که مدیران باید کاری را که به دیگران توصیه می‌کنند ابتدا خود آن را انجام دهند

3- به هنگام برقراری ارتباط با مشتری، از نوآوری و انجام ریسک‌های حساب شده استقبال کنید. زمانی که کارکنان سازمان برای انجام کارهای نو و بدیع در پاسخ به مشتریان احساس آزادی کنند، فرهنگ ارتباط با مشتری بیشتر تقویت می‌شود

4- برای حفظ ارتباط با مشتری مسوولیت‌های هر فرد را روشن کنید. مدیران باید راهنمایی و هدایت کافی برای هماهنگ سازی فعالیت‌های مختلف را فراهم کنند و مطمئن شوند که اطلاعات مشتری به افراد مناسبی منتقل می‌شود

5- افراد را به مشارکت در تصمیم‌گیری‌های مربوط به ارتباط با مشتریان تشویق کنید. مدیران باید به تمام کسانی که به اطلاعات مشتریان دسترسی دارند، اختیار و قدرت تصمیم‌گیری بدهد. علاوه بر این اقدامات، مدیرانی که در ایجاد فرهنگ ارتباط با مشتری موفق عمل کردند گرایش‌ها و جهت‌گیری‌های دیگری نیز داشته اند. آنها در کنار توجه به مشتری روی رقبا، عوامل موثر در بازار و جامعه متمرکز شده اند. سعی آنها بر این است که از بوروکراسی، سیاست‌های اداری و هر آنچه مخمل عملکرد سازمان می‌شود اجتناب ورزند.

عامل اصلی در شکل‌گیری فرهنگ ایجاد ارتباط با مشتری اندازه یا پیچیدگی سازمان‌ها نیست بلکه خواست و مهارت مدیران ارشد نقش اصلی را ایفا می‌کند

۴-۵ مزایای بکارگیری سیستم‌های مدیریت ارتباط با مشتری

CRM یک واقعیت ملموس برای سازمان‌های تجاری است و به طور خلاصه مزایای زیر را برای سازمان به دنبال دارد:

- ۱- پاسخگویی سریع به درخواست مشتریان
- ۲- فراهم کردن شرایط مراجعه مجدد مشتری
- ۳- کاهش هزینه‌های تبلیغاتی
- ۴- افزایش فرصت‌های بازاریابی و فروش
- ۵- شناخت عمیق تر مشتری
- ۶- دریافت باز خورد از مشتری و توسعه خدمات و محصولات جاری

۴-۶ چارچوب گارتنر

طرح‌های CRM نیازمند چارچوبی اند که تضمین کند برنامه‌های سازمان در مبنای استراتژیک و یکپارچه در نظر گرفته می‌شوند. گارتنر چنین رویکردی را که شامل 8 گام است طراحی کرده است:

1- تدوین چشم انداز سازمان ایجاد چشم انداز موثر مستلزم این است که رهبران سازمانی: • معانی CRM را برای موسسه تعریف کنند • اهداف را تعیین کنند • تصویری از آنچه سازمان می‌خواهد برای مشتریان هدفش باشد ترسیم کنند

هدف: باید مجموعه‌ای از ارزش‌های متمایز شده که برای مشتریان مهم است خلق گردد.

2- تدوین استراتژیهای CRM استراتژی CRM نگرشی را در مورد نحوه ایجاد ارتباط با مشتریان ارزشمند و نحوه وفاداری در آنان فراهم می‌آورد. گام اول: تدوین استراتژی CRM بخش بندی مشتریان در گروه‌ها، تعیین اهداف و معیارهای سنجش برای هر بخش است. گام دوم: ارزیابی وضعیت پایگاه مشتری به عنوان یک دارایی است. این کار از طریق ترسیم نمودار قوت‌ها و ارزش ارتباط با مشتریان در دو بعد صورت می‌گیرد:

• مشتری تا چه اندازه برای سازمان ارزشمند است؟

• سازمان تا چه اندازه برای مشتری ارزشمند است؟

3- طراحی تجربه مشتری در این مرحله باید اطمینان حاصل شود که محصولات و تعاملات سازمان باعث خلق ارزش برای مشتریان گشته، به طور پایدار ارائه شده و به موقعیت به بازار مطلوب دست پیدا کرده است یا خیر اجرای سیستم باز خورد عملیاتی سبب افزایش آگاهی سازمان از شکایات مشتری گشته، حل این شکایات را میسر می‌سازد.

4- میسر ساختن همکاری سازمانی همکاری سازمانی به معنی تغییر فرهنگ، ساختارهای سازمانی و رفتارهاست تا اطمینان حاصل شود که کارکنان، شرکا و تامین کنندگان در جهت ایجاد ارزش برای مشتریان با یکدیگر همکاری می‌کنند.

5- طراحی مجدد فرآیندهای کسب و کار در زمینه طراحی مجدد فرآیندهای مربوط به مشتری استفاده از چارچوب زیر الزامی است: • نقاط تماس و فرآیندهایی را که بر روی مشتریان تاثیر می‌گذارند، حسابرسی کرده، نقشه آنها را ترسیم کنید. • فرآیندهای کلیدی را از دیدگاه مشتری شناسایی کنید و فرآیندهایی که می‌توانند بیشترین نارضایتی را ایجاد کنند پیدا کرده، در حله اول بر آنها تمرکز شود.

- این فرآیندها را بر اساس اثراشان بر روی هدف CRM مشخص و اولویت بندی کنید.
 - تغییرات لازم را در سازمان اعمال کنید (هیچ فرآیندی را نباید بدون مجری و مسئول رها کرد)
 - با استفاده از اهداف مشتریان، اهداف قابل اندازه گیری و یا مفهومی را تعیین کنید. برای هر فرآیند کلیدی نوعی توافقنامه سطح خدمت به مشتری را برقرار کنید.
 - فرآیندها را بر اساس اهمیت آنها به مشتریان و اثرات آنها بر اهداف CRM الویت بندی کنید.
- 6- تدوین استراتژی اطلاعات مشتری منظور از این مرحله جمع آوری داده‌های صحیح و ارسال آنها به مکان صحیح است. مدیریت موفق ارتباط با مشتری نیازمند خلق نوعی ((عرضه خون)) اطلاعاتی است که در سراسر سازمان جریان یافته و سیستم‌های عملیاتی و تحلیلی را یکپارچه کند.
- 7- استفاده از فناوری منظور از این مرحله مدیریت داده‌ها و اطلاعات، برنامه‌های کاربردی پیش روی مشتری، زیر ساخت‌ها و معماری IT است.
- 8- معیارهای سنجش منظور از این مرحله اندازه گیری شاخص‌های درونی و بیرونی موفقیت و شکست CRM است.

این شاخص‌ها دارای کاربرد زیر هستند:

- تعیین و اندازه گیری سطح تحقق اهداف CRM
 - ارائه باز خورد برای اصلاح استراتژی CRM و اجرای آن
 - نظارت بر تجربه مشتری از سازمان
 - تغییر شیوه جبران خدمات کارکنان و مشوق‌های داده شده
 - ارزیابی سازمان نسبت به رقبا
- در گام آخر: اهدافی را که باید برآورده شوند و تاکتیک‌هایی را که باید مورد استفاده قرار گیرند تعریف می‌کنیم.

۴-۷ اندازه گیری رضایت مشتری (CSM) چیست؟

اندازه گیری رضایت مشتریان یکی از فعالیت های بسیار مهم برای هر سازمان است. سازمانی که بتواند رضایت مشتریان خودش را اندازه گیری نماید، در حقیقت مانند این است که علائم حیاتی خودش را دائماً پیش رو دارد، و می تواند آنها را کنترل کند و با توجه به آنها سازمان خودش را به پیش برد.

مروزه در دنیای مدرن و مدیریت مدرن اندازه گیری رضایت مشتریان از پیشرفت الگوی مدرن تبدیل به یک ضرورت شده است

اهمیت CSM در استاندارد بند 1-2-8- استاندارد ISO 9001 بیان می دارد :

به عنوان یکی از طرق اندازه گیری عملکرد سیستم مدیریت کیفیت، سازمان بایست اطلاعات مرتبط با برداشت و درک مشتریان را پالایش نموده و تعیین نمایند که آیا سازمان توانسته نیازمندی های مشتریان را برآورده سازد. ضمناً "متمدهای دستیابی و به کارگیری این اطلاعات نیز بایستی تعیین گردند .

4-7-1 نحوه اندازه گیری رضایت مشتری چگونه است؟

- 1- نمونه گیری: (در نظر گرفتن تمام جوانب) پراکندگی جغرافیایی و پراکندگی از لحاظ بازار
- 2- ابزار تحقیق: تعیین شاخص های رضایت که مهم ترین عناصر در رضایت مشتریان را اندازه گیری کند .
- 3- متدولوژی تحقیق: انتخاب روش آماری و تحلیل با هدف آسان شدن عددی کردن شاخص های کیفی
- 4- تشخیص منابع: تعریف منابع و استراتژی سازمان پس از تحلیل در فعالیتهای شرکت

اثرات رضایت مشتری چیست؟

رضایت مشتری از سه طریق به افزایش درآمد و سود منجر می شود. تکرار خرید مشتری، خرید کالای جدید و خرید کالا توسط مشتریان جدیدی که توسط مشتریان راضی به کالا تمایل پیدا کرده اند. امروزه در کشورهای صنعتی برنامه های ارتباط با مشتریان در سرلوحه برنامه های

بازاریابی عرضه کنندگان قرار گرفته است. امروزه دیگر هیچ تولیدکننده و عرضه کننده ای به فروش یک بار به مشتری نمی اندیشد .

نکته مهم این است که در صورت نارضایتی مشتری، تمامی مکانیسم ها در جهت عکس عمل کرده و درآمد و سود عرضه کننده را کاهش می دهد. بررسیها نشان می دهد که مشتریان ناراضی در انتقال احساس خود به دیگران فعال تر و موفق تر عمل می کنند .

۴-۸ تفاوت CRM و CSM چیست؟

هدف یک سیستم مدیریت روابط با (مشتریان، داشتن مشتریان وفادار است). یعنی شما با اجرای CRM در سازمانتان انتظار دارید مشتریانی راضی تر (حتی اگر کمتر از قبل هستند) داشته باشید تا پول بیشتری از مشتریان کمتر بدست بیارید. سپس حتما می خواهید تعداد آنها را افزایش بدهید. این در درازمدت به شما کمک میکند که مشتریانی راضی و بسیار داشته باشید .

ولی CSM در ارتباط با میزان رضایتمندی مشتریان شما، صرفا یکی از عوامل (البته از عوامل مهم CRM هست). هر مشتری به دنبال یک نیاز به دنبال محصولی میرود و بعد از یافتن خواسته هایش در محصول اقدام به خرید می کند. ولی وفاداری مشتری که هدف سیستم های CRM هست فقط از طریق جلب رضایت آنها به دست می آید.

فصل پنجم

ارتباط بین داده

کاوی و مدیریت

ارتباط با مشتری

۵-۱ ارتباط بین داده کاوی و مدیریت ارتباط با مشتری

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

برای روشن تر شدن مساله می توان مثال را این گونه بیان کرد که در یک فروشگاه زنجیره ای پس از داده کاوی مشخص می شود که درصدی از مشتریان خرید تلویزیون ، میز تلویزیون و گلدان کریستالی را هم در همان روز و بعد از خرید تلویزیون می خرند. مدیر فروشگاه می تواند بلافاصله دستوراتی صادر کند که بر اساس مدل های تلویزیون موجود میزهایی و براساس مدل میزها گلدان های کریستالی برایفروش سفارش داده شود و غرفه های جنبی غرفه تلویزیون را به میز و گلدان کریستالی اختصاص دهد . مطمئنا" 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 یک تصویر از واقعیت به شکل مدل می سازد که این مدل روابط موجود در داده ها را شرح می دهد. برای بهبودی بهره وری از یک فروشگاه داده کاوی از داده های انبار داده، مدل هایی را ارائه می دهد که بیانگر این هستند که چه محصولاتی یا خدماتی ،به چه مشتریانی،در چه زمانی واز طریق چه کانالی عرضه شود.

بیشتر شرکت ها بانک های 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 ارتباط با مشتری است. رابطه مشتری با زمان تغییر می کند و چنانچه تجارت و مشتری درباره یکدیگر بیشتر بدانند این رابطه تکامل ورشد 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 مشتری می گوید چه چیزی احتمالا" جالب توجه است و چه تصمیماتی باید گرفته شود. داده کاوی می تواند سود آوری مشتری های بالقوه را که می توانند به مشتریان بالفعل تبدیل شوند،پیش بینی کند و اینکه مشتری تا چه زمانی وفادار خواهد ماند وچگونه احتمالا" مارا ترک خواهد کرد.بعضی از مشتریان مرتبا" مراجعاتشان رابه شرکت ها برای کسب مزیت هایی که طی رقابت میان آنها به وجود می آید، تغییر می دهند.

در این صورت شرکت ها می توانند هدفشان را روی مشتریانی متمرکز کنند که سود آوری 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.

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 و مراجع

- 1- کتاب سال شبکه (کتاب مرجع صنعت شبکه در ایران) اذر ماه سال 1390 مقاله دوم peer1 hosting
- 2- کتاب مدیریت کیفیت فراگیر (جلد اول) دکتر مصطفی جعفری موسسه خدمات فرهنگی رسا
- 3- مقاله داده کاوی، مهندس مهرداد اشکانندی راد
- 4- Fayyad U., Piatetsky-Shapiro G., and Smyth P., "From Data Mining to Knowledge Discovery in Databases," American Association for Artificial Intelligence, 1996.
- 5- Introduction to Data Mining and Knowledge Discovery By Two Crows Corporation
- 6- Top 10 algorithms in data mining XindongWu Knowl Inf Syst (2008)
- 7- سایت <http://www.wikipedia.org>
- 8- مقاله مدیریت ارتباط با مشتری، پرستو شاه محمدی
- 9- سایت <http://irandataminer.ir> در مورد وب کاوی
- 10- پایان نامه خانم شیما شاه حسینی، دانشگاه علم و صنعت، وب کاوی، 1389
- 11- سایت www.modiriran.ir در رابطه با مدیریت ارتباط با مشتری
- 12- مقاله آقای ارمان محمد زاده، دانشگاه ازاد اسلامی واحد خمینی شهر، وب کاوی، 1391