

# آشنایی با مفاهیم و روش‌های داده‌کاوی



# آشنایی با مفاهیم و روش‌های داده‌کاوی

عباس مرادی  
جواد حسین زاده  
اشکان شباک  
کاوه کیانی  
محمد شیرینی



پژوهشکده‌ی انبار



بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِیْمِ

## پیش‌گفتار

در این کتاب روش‌های داده‌کاوی تشریح می‌شود. داده‌کاوی به معنای کشف دانش و استخراج از مقادیر زیادی از داده‌های خام است. از آنجا که تنها ابزار مفید برای استفاده از تحلیل وقایع گذشته در پایگاه داده‌ها به منظور پیش‌بینی در آینده، داده‌کاوی خواهد بود. از این رو برای کشف و استخراج اطلاعات از آمار مد نظر برنامه‌ریزان و مدیران هر سازمانی نیازمند داده‌کاوی خواهیم بود. از سوی دیگر عصر حاضر، عصر روش‌های نوین و متنوع تولید، ذخیره و گردآوری داده‌ها و اطلاعات آماری است. سیاست‌گذاران با حجم انبوهی از اطلاعات آماری متنوع که در کسری از زمان بیشتر و بیشتر تولید می‌شوند در سازمان‌های مهم و کلیدی کشور مواجه هستند. داده‌کاوی اسلوبی برای یافتن دانش و اطلاعات بهینه از میان این انبوه از ریزداده‌ها است. دسترسی به الگوهای مخفی داده‌ها و اطلاعات آماری انباشت شده توسط سازمان‌ها حاکی از آن است که داده‌کاوی بسیار ارزشمند و گران‌بهاست. کتاب حاضر حاوی مفاهیم و روش‌های اساسی داده‌کاوی است و در ۵ فصل تنظیم شده است. در انتهای هر فصل پس از بیان خلاصه به توضیح منابع و مراجع برای مطالعه‌ی بیشتر نیز پرداخته خواهد شد. فصل اول به مقدمه‌ای در مورد

مفاهیم داده‌کاوی می‌باشد. این فصل مشتمل بر چشم‌انداز و مأموریت کتاب برای آشنایی با سایر مطالب کتاب خواهد بود. فصل دوم و سوم به مفاهیم داده‌کاوی و آماده‌سازی (پاک‌سازی و تمیزسازی) داده‌ها اختصاص دارد. در فصل چهارم به پردازش‌های برخط، انبار داده‌ها و پایگاه داده پرداخته خواهد شد. در فصل پنجم روش‌های داده‌کاوی شامل درخت تصمیم، خوشه‌بندی، رده‌بندی، قواعد پیوند و سایر الگوریتم‌های داده‌کاوی مورد مطالعه قرار می‌گیرد. فصل آخر به برخی از موضوعات مرتبط با داده‌کاوی پرداخته خواهد شد.

---

# فهرست مطالب

مقدمه	۱
۱-۱ مقدمه	۱
۲-۱ تعاریف و مفاهیم داده‌کاوی	۴
کلیات داده‌کاوی	۲۳
۱-۲ مفاهیم بنیادی داده‌کاوی	۲۳
۲-۲ تحلیل‌های پیش‌بینی‌کننده و تحلیل‌های توصیفی	۲۵
۳-۲ ضرورت استفاده از داده‌کاوی	۳۱
۴-۲ برخی از کاربردهای مهم داده‌کاوی	۳۲
آماده‌سازی داده‌ها	۳۵
۱-۳ آماده‌سازی داده‌ها	۳۵
۲-۳ تبدیلات عددی	۳۹
۳-۳ گروه‌بندی کردن	۴۱
۴-۳ انبوهش و خلاصه‌سازی داده‌ها	۴۱
۵-۳ داده‌های گم‌شده	۴۱
۶-۳ حذف داده‌های پرت	۴۲
۷-۳ انتخاب زیرمجموعه‌ایی از ویژگی‌ها و پالایش نمونه‌ها	۴۳
۸-۳ خلق ویژگی	۴۳
پردازش‌های برخط	۴۵
۱-۴ مقدمه	۴۵
۲-۴ پردازش‌های تحلیلی برخط (OLAP)	۴۵
۳-۴ احداث مکعب‌های OLAP از انبارداده‌ها	۵۲
۴-۴ پردازش‌های تراکنشی برخط (OLTP)	۵۳
روش‌های داده‌کاوی	۴۵

۵-۱	مقدمه	۵۵
۵-۲	الگوریتم درخت تصمیم	۵۵
۵-۴	الگوریتم خوشه‌بندی	۹۴
۵-۵	الگوریتم رگرسیون	۱۱۴
۵-۶	الگوریتم بیزساده	۱۲۱
۵-۸	الگوریتم سری زمانی	۱۳۱
۵-۹	الگوریتم شبکه‌های عصبی	۱۴۵
۵-۱۰	الگوریتم قواعد پیوند	۱۶۰
۵-۱۱	انواع تقسیم‌بندی‌های داده‌کاوی	۱۶۳
۵-۱۲	برخی از زبان‌های برنامه‌نویسی داده‌کاوی	۱۶۴
	<b>داده‌کاوی در آمار رسمی</b>	۱۷۲
۶-۱	مقدمه	۱۷۲
۶-۲	آمار رسمی	۱۷۲
۶-۳	اصول بنیادی آمارهای رسمی	۱۷۳
۶-۴	بررسی الگوی مصرف خانوارهای شهری براساس طرح هزینه و درآمد خانوار ۱۷۶	
۶-۵	کاربرد داده‌کاوی در پزشکی - مطالعه موردی تشخیص دیابت با استفاده از چربی خون (استفاده داده‌کاوی از آمار ثبتي در آزمایشگاه‌ها)	۱۸۶
۶-۶	نقش داده‌کاوی در «چارچوب کلی فرایند کسب و کار آماری»	۱۸۶
	<b>پیوست</b>	۱۹۷
	<b>مرجع</b>	۲۰۱



# فهرست جداول

- جدول ۱-۳: گسسته‌سازی متغیر پیوسته‌ی سن در مرکز آمار ایران ..... ۴۰
- جدول ۱-۴: وضعیت اشتغال افراد و ویژگی‌های آن‌ها ..... ۴۶
- جدول ۱-۵: نمونه‌ای از داده‌های آموزشی ..... ۶۵
- جدول ۲-۵: سه مشاهده با چهار متغیر مفروض ..... ۱۰۴
- جدول ۳-۵: تفسیر ضریب نیم‌رخ ..... ۱۱۳
- جدول ۴-۵: انواع مدل‌های رگرسیون ..... ۱۲۰



# فهرست شکل

- شکل ۱-۱: معماری انبارداده ..... ۱۰
- شکل ۲-۱: شمای انبارداده ستاره‌ای ..... ۱۹
- شکل ۳-۱: شمای انبارداده برفگونه ..... ۱۹
- شکل ۴-۱: پردازش‌های داده‌کاوی در یک نگاه ..... ۲۲
- شکل ۱-۲: هرم دانش ..... ۲۵
- شکل ۲-۲: ساختار مغز انسان ..... ۲۸
- شکل ۳-۲: نمایی کلی از ساختار شبکه‌های عصبی ..... ۲۸
- شکل ۱-۳: بیست نرم‌افزار برتر دیدارسازی داده‌ها تا سال ۲۰۱۹ میلادی ..... ۳۷
- شکل ۱-۴: بعدهای وضعیت اشتغال ..... ۴۷
- شکل ۲-۴: مکعب وضعیت اشتغال ..... ۴۸
- شکل ۳-۴: سلسله مراتب زمان ..... ۴۹
- شکل ۴-۴: سلسله مراتب مکان ..... ۴۹
- شکل ۵-۴: ایجاد یک زیرمکعب از مکعب اصلی ..... ۵۱
- شکل ۶-۴: خرد کردن یک مکعب به مکعبی با مختصات بیشتر (تظریف یک مکعب) ..... ۵۲
- شکل ۱-۵: داده‌های دانش‌آموزان شهرستان در یک نگاه ..... ۵۶
- شکل ۲-۵: درخت تصمیم ..... ۵۷
- شکل ۳-۵: تاثیرات متغیرهای مهم در ادامه تحصیل ..... ۵۸

- شکل ۴-۵ : تاثیرگذاری متغیرهای ورودی بر متغیر هدف ..... ۵۸
- شکل ۵-۵ : بیشترین تاثیرگذاری متغیرهای ورودی بر متغیر هدف ..... ۵۹
- شکل ۷-۵ : درخت تصمیم در حالت کلی برای داده‌های آموزشی مثال فوق ..... ۶۷
- شکل ۸-۵ : درخت تصمیم با جزئیات ..... ۶۸
- شکل ۹-۵ : کاهش بهره‌ی اطلاعات از ریشه به برگ‌ها ..... ۷۲
- شکل ۱۰-۵ : رده‌بندی ..... ۷۴
- شکل ۱۱-۵ : رده‌بندی با استفاده از درخت تصمیم ..... ۷۴
- شکل ۱۲-۵ : انواع جداکننده‌ها ..... ۷۷
- شکل ۱۳-۵ : در جستجوی بهترین جداساز ..... ۷۷
- شکل ۱۴-۵ : بهترین جداساز (SVM) ..... ۷۸
- شکل ۱۵-۵ : داده‌های مورد مطالعه ..... ۸۴
- شکل ۱۶-۵ : مشخص کردن ویژگی مورد نظر ..... ۸۴
- شکل ۱۷-۵ : یافتن رده‌ایی خاص از داده‌ها ..... ۸۵
- شکل ۱۸-۵ : جمعیتی از کروموزوم‌ها ..... ۸۸
- شکل ۱۹-۵ : فلوجارت الگوریتم ژنتیک ..... ۹۰
- شکل ۲۰-۵ : رده‌بندی نرم (تغییر رنگ‌ها به آرامی صورت می‌گیرد) ..... ۹۳
- شکل ۲۱-۵ : رده‌بندی سخت (تغییر رنگ‌ها سریع صورت می‌گیرد و دو رنگ داریم) ..... ۹۳
- شکل ۲۲-۵ : برخی از انواع خوشه‌بندی از نظر ظاهری ..... ۹۸
- شکل ۲۳-۵ : نمودار درختی خوشه‌بندی سلسله‌مراتبی تجمعی و تقسیمی ..... ۱۰۱
- شکل ۲۴-۵ : نمودار دندروگرام خوشه‌بندی سلسله‌مراتبی ..... ۱۰۲
- شکل ۲۵-۵ : خوشه‌بندی با استفاده از روش k-Means (نتایج نادرست) ..... ۱۰۷
- شکل ۲۶-۵ : خوشه‌بندی با استفاده از روش مبتنی بر چگالی (نتایج درست) ..... ۱۰۷
- شکل ۲۷-۵ : خوشه‌بندی نرم ..... ۱۰۸
- شکل ۲۸-۵ : توزیع مشاهده‌ها (یک بعدی) ..... ۱۰۹
- شکل ۲۹-۵ : خوشه‌بندی سخت (تفکیکی- مبتنی بر افراز) ..... ۱۰۹
- شکل ۳۰-۵ : خوشه‌بندی نرم ..... ۱۰۹

- شکل ۵-۳۱ : رگرسیون خطی ..... ۱۱۵
- شکل ۵-۳۲ : رگرسیون سهمی ..... ۱۱۶
- شکل ۵-۳۳ : رگرسیون درجه ۳ ..... ۱۱۷
- شکل ۵-۳۴ : مدل مفهومی رگرسیون لجستیک ..... ۱۱۸
- شکل ۵-۳۵ : فضای نمونه‌ی افزاز شده به  $k$  کلاس ( $k$  رده) ..... ۱۲۱
- شکل ۵-۳۷ : جداول رابطه‌ی مشتریان و سیاهه‌ی کلیک آن‌ها ..... ۱۲۷
- شکل ۵-۳۸ : روند کاهشی ..... ۱۳۵
- شکل ۵-۳۹ : روند ثابت ..... ۱۳۵
- شکل ۵-۴۰ : روند افزایشی ..... ۱۳۶
- شکل ۵-۴۱ : تناوب ..... ۱۳۷
- شکل ۵-۴۲ : تغییرات فصلی ..... ۱۳۸
- شکل ۵-۴۳ : تغییرات غیرمعمولی ..... ۱۳۹
- شکل ۵-۴۴ : میانگین متحرک مرتبه چهار و هفت برای یک سری زمانی واقعی ..... ۱۴۲
- شکل ۵-۴۵ : درخت اتورگرسیو (AR Tree) ..... ۱۴۴
- شکل ۵-۴۶ : تفاوت AR و ART ..... ۱۴۴
- شکل ۵-۴۷ : ساختار مغز انسان ..... ۱۴۵
- شکل ۵-۴۸ : نمایی کلی از ساختار شبکه‌های عصبی- با دولایه‌ی پنهان ..... ۱۴۷
- شکل ۵-۴۹ : مدل کلی شبکه‌های عصبی ..... ۱۴۹
- شکل ۵-۵۰ : خروجی شبکه‌ی عصبی ..... ۱۵۰
- شکل ۵-۵۱ : یک واحد پردازش پایه‌ای در شبکه‌های عصبی ..... ۱۵۲
- شکل ۵-۵۲ : معماری شبکه عصبی پیش‌رو ..... ۱۵۶
- شکل ۵-۵۳ : معماری شبکه عصبی بازگشتی ..... ۱۵۶
- شکل ۵-۵۴ : داده‌های فرضی شامل مربع و جمع - جداسازی غیر خطی ..... ۱۵۸
- شکل ۵-۵۵ : محاسبه‌ی معیارهای قواعد پیوند ..... ۱۶۱
- شکل ۵-۵۶ : تقسیم‌بندی برخی از الگوریتم‌های داده‌کاوی ..... ۱۶۴
- شکل ۵-۵۷ : شش زبان برتر برنامه‌نویسی علوم داده در سال ۲۰۱۹ ..... ۱۶۵

- شکل ۵-۵۸: آخرین ویژگی‌های پایتون ..... ۱۶۶
- شکل ۶-۱: پایگاه داده رابطه‌ی هزینه و درآمد خانوار شهری سال ۱۳۹۰ ..... ۱۷۸
- شکل ۶-۲: پایگاه داده رابطه‌ی هزینه و درآمد خانوار شهری سال ۱۳۹۶ ..... ۱۷۹

# فصل ۱

## مقدمه

### ۱-۱ مقدمه

داده‌کاوی به معنای کشف دانش و استخراج از مقادیر زیادی از داده‌های خام است. در واقع داده‌کاوی مفیدترین ابزار برای استفاده از تحلیل وقایع گذشته در پایگاه‌های داده‌ای خواهد بود. از این رو برای کشف و استخراج اطلاعات از داده‌ها و اطلاعات آماری یک سازمان نیازمند استفاده از داده‌کاوی خواهیم بود. با این اطلاعات کشف شده تصمیم‌گیری‌ها و برنامه‌ریزی‌ها بسیار دقیق‌تر و کاربردی‌تر خواهد شد. سیاست‌گذاران نیازمند روش‌های داده‌کاوی برای برنامه‌ریزی و سیاست‌گذاری‌های کلان خواهند بود. داده‌کاوی بین محققان علوم آمار، کامپیوتر و ریاضی در جهت تحلیل داده‌ها و پایگاه‌های داده‌ای که شامل مه‌داده‌ها و یا داده‌های با ابعاد بالا هستند، به عنوان یکی از ابزارهای پرتوان شناخته شده است. بر همین اساس مرکز آمار ایران به عنوان یگانه مرجع آمار رسمی کشور روز به روز به اهمیت داده‌کاوی پی برده و این مهم را در اولویت‌های برنامه‌ای خود قرار داده است. پرواضح است که ایجاد یک نظام کارآمد و مؤثر در تولید و عرضه

آمار از الزامات اولیه و ضروری در برنامه‌ریزی است و زیربنای برنامه‌ریزی مناسب، اطلاعات جامع، منسجم و به‌روز است. داده‌کاوی به معنای کشف دانش و استخراج از مقادیر زیادی از داده‌های خام است. از این رو داده‌کاوی مفیدترین ابزار برای استفاده از تحلیل وقایع گذشته در پایگاه‌های داده‌ای مراکز آمار ایران خواهد بود. بدین معنی برای کشف و استخراج اطلاعات از آمار رسمی کشور که آمار مد نظر، سیاست‌گذاران، برنامه‌ریزان و مدیران کشور است نیازمند استفاده از داده‌کاوی خواهیم بود. داده‌کاوی دانشی بین‌رشته‌ای است و ترکیبی از علوم مانند آمار، ریاضیات، کامپیوتر، علوم اطلاعات، نظریه پایگاه داده و یادگیری ماشین می‌باشد. از آنجا که داده‌کاوی یک علم بین‌رشته‌ای است، فلسفه پیدایش آن در درازمدت شکل گرفته است. به طوریکه ابتدا در قرن هجدهم قضیه بیز توسط توماس بیز چاپ و منتشر گردید که بعدها احتمال شرطی نام گرفت. این قضیه از آن جهت مفید است که می‌توان از طریق آن احتمال یک پیشامد را با مشروط کردن نسبت به وقوع و یا عدم وقوع یک پیشامد دیگر محاسبه کرد. این دیدگاه پایه و اساس داده‌کاوی و احتمالات است. قرن بعد [قرن نوزدهم]، قرن پیدایش الگوریتم رگرسیون توسط گاوس بود. ایشان از رگرسیون برای تعیین چرخش اجرام به دور خورشید استفاده کردند. هدف از تحلیل‌های رگرسیونی، یافتن تخمین و روابط مناسب بین متغیرها است. رگرسیون شامل تکنیک‌های زیادی برای مدل سازی و تحلیل متغیرهای خاص و منحصر به فرد است. مادامی که هدف یافتن روابط بین متغیر وابسته و یک یا چند متغیر مستقل باشد، می‌توان از این الگوریتم به درستی استفاده کرد. تحلیل رگرسیون به صورت گسترده برای پیش‌بینی استفاده شده و یکی از ابزارهای کلیدی در داده‌کاوی به حساب می‌آید. به کمک رگرسیون پیش‌بینی‌هایی از داده‌های خام برای امور آتی صورت می‌گیرد. در قرن بعد [قرن بیستم]، الگوریتم شبکه‌های عصبی توسط دانشمندان کشف و منتشر شد. این ایده برگرفته از شبکه‌های عصبی مغز بود. در سال ۱۹۴۳ دو پژوهش‌گر به نام‌های وارن مکلوج و پیتس مقاله‌ای در مورد شبکه‌های عصبی نوشتند [۱۱]. این مقاله اولین و جرقه در زمینه‌ی طراحی شبکه‌های عصبی مصنوعی، برای سایر دانشمندان بود. پس از این مقاله انواع مختلفی از شبکه‌های عصبی مصنوعی توسط اندیشمندان تهیه و تنظیم شد، که همگی برگرفته از ایده‌های این دو دانشمند شهیر بود. تمامی این شبکه‌های



مصنوعی متفاوت و گوناگون دارای یک وجه مشترک بودند که همگی از مغز انسان اقتباس لازم را گرفته بودند و نودها و گره‌ها را بر اساس نورون‌های مغز می‌ساختند و یال‌ها و ارتباطات بین گره‌ها را بر اساس آکسون‌های مغز می‌ساختند. حدود بیست سال بعد یعنی در سال ۱۹۶۶ میلادی شرکت جدیدی به نام علوم تصمیم‌گیری توسط لورنس فوگل تأسیس گردید. این شرکت در واقع اولین شرکتی بود که به طور ویژه از محاسبات تکاملی برای حل مسائل واقعی استفاده می‌کرد [۷]. در سال ۱۹۷۰ سیستم‌های مدیریت پایگاه داده پیشرفته اختراع گردید و امکان ذخیره‌سازی داده‌ها برای همگان فراهم گشت. پنج سال بعد یعنی سال ۱۹۷۵ میلادی، جان هنری هلند کتابی به نام الگوریتم ژنتیک نوشت این کتاب تکنیک جستجویی در علم رایانه برای یافتن راه‌حل تقریبی برای بهینه‌سازی و مسائل جستجو است [۹]. در واقع الگوریتم‌های ژنتیک از اصول انتخاب طبیعی داروین برای یافتن فرمول بهینه جهت پیش‌بینی یا تطبیق الگو استفاده می‌کنند. الگوریتم‌های ژنتیک اغلب گزینه خوبی برای تکنیک‌های پیش‌بینی بر مبنای رگرسیون هستند. در هوش مصنوعی الگوریتم ژنتیک (یا GA) یک تکنیک برنامه‌نویسی است که از تکامل ژنتیکی به عنوان یک الگوی حل مسئله استفاده می‌کند. اما در نهایت پنج سال بعد یعنی سال ۱۹۸۰ میلادی، برای اولین بار اصطلاح «داده‌کاوی» مطرح شد. در این دوره کارشناسان می‌توانستند به روابط معنادار پی ببرند. سپس، در سال ۱۹۸۹ اصطلاح «کشف دانش در پایگاه داده» مطرح شد. در همین زمان اولین کارگاه آموزشی با نام KDD برای استخراج دانش در پایگاه‌های داده گوناگون، شروع به فعالیت نمود. در سال ۲۰۰۱ میلادی، علم داده به عنوان یک رشته مستقل معرفی شد. تا اینکه در قرن حاضر داده‌کاوی به معنای واقعی خود نمایان شد. امروزه داده‌کاوی در اقتصاد، مهندسی، سیاست، پزشکی و... کاربرد فراوانی دارد. داده‌کاوی تراکنش‌های مالی، سود سهام، شبکه‌های عصبی، امنیت ملی، الگوریتم‌های ژنتیک نمونه‌هایی از این کاربردهاست. در سال ۲۰۱۴ میلادی آقای حسنی و همکاران به بررسی اهمیت و ارزش کاربردهای داده‌کاوی در آمار رسمی پرداختند [۸].

## ۲-۱ تعاریف و مفاهیم داده‌کاوی

به منظور برداشت یکسان از واژه‌ها و اصطلاحات به کار رفته در این کتاب، تعاریف و مفاهیم برخی واژه‌های کلیدی در زیر آرایه می‌شود.

**داده<sup>۱</sup>:** صفت یا اطلاعی است، اغلب عددی، که از طریق مشاهده گردآوری می‌شود. داده‌ها به دو صورت داده‌های خرد و کلان تقسیم‌بندی می‌شوند. در حقیقت، داده‌ها حقایق و واقعیت‌های خامی هستند که اطلاعات از آن‌ها ساخته می‌شوند. این اجزاء در پایگاه‌های داده، ذخیره و مدیریت می‌شوند.

**اطلاعات<sup>۲</sup>:** اضافه کردن زمینه و تفسیر به داده‌ها و ارتباط آن‌ها به یکدیگر، موجب شکل‌گیری اطلاعات می‌شود. اطلاعات در حقیقت داده‌های خلاصه شده را در برمی‌گیرد که گروه‌بندی، ذخیره، پالایش، سازماندهی و تحلیل و تفسیر شده‌اند.

**اطلاعات آماری<sup>۳</sup>:** منظور از اطلاعات آماری، نوع خاصی از اطلاعات هستند که با انجام عملیات‌های ریاضی بر روی داده‌ها و یا دسته‌بندی و طبقه‌بندی آن‌ها بر اساس یک روش تعریف شده به دست آمده‌اند. از این مفهوم برای توصیف بهتر تعریف آمار از بعد عددی آن استفاده می‌شود.

**فراداده‌ها<sup>۴</sup>:** اطلاعاتی درباره‌ی داده‌ها و فرایندهای تولید و نحوه استفاده از داده‌ها است. **پایگاه داده<sup>۵</sup>:** مجموعه‌ای ساختاریافته از داده‌هاست که در یک سیستم رایانه‌ای نگهداری می‌شود و از راه‌های مختلفی در دسترس قرار می‌گیرد.

**دانش<sup>۶</sup>:** دانش مجموعه‌ای از شناخت‌ها و مهارت‌های لازم برای حل یک مسأله است لذا اگر اطلاعاتی که در دست است بتواند مشکلی را حل کند می‌توان گفت دانش وجود دارد. ضمن آنکه دانش باید امکان تبدیل به دستورالعمل اجرایی و عملی شدن را داشته باشد.

**مدیریت دانش<sup>۷</sup>:** به معنای در دسترس قرار دادن نظام‌مند اطلاعات و اندوخته‌های علمی است، به گونه‌ای که به هنگام نیاز در اختیار افرادی که نیازمند آن‌ها هستند، قرار گیرند تا

<sup>1</sup> Data

<sup>2</sup> Information

<sup>3</sup> Statistical Information

<sup>4</sup> Metadata

<sup>5</sup> Database

<sup>6</sup> Knowledge

آن‌ها بتوانند کار روزمره خود را با بازدهی بیشتر و مؤثرتر انجام دهند. مدیریت دانش شامل یک سری استراتژی و راهکار برای شناسایی، ایجاد، نمایندگی، پخش و تطبیق بینش‌ها و تجارب در سازمان می‌باشد.

**داده‌کاوی<sup>۸</sup>:** روش‌های کشف الگوها در مجموعه داده‌های بزرگ داده‌کاوی نام دارد. داده‌کاوی دانش و مهارتی است که بین علوم مختلفی همچون یادگیری ماشین، آمار و کامپیوتر قرار دارد که هدف کلی آن استخراج اطلاعات با روش‌های هوشمند از مجموعه‌ای از داده‌هاست به نحوی که اطلاعات کشف شده به وضوح ملموس‌تر از گذشته باشد.

**آمار رسمی<sup>۹</sup>:** به اطلاعات عددی گفته می‌شود که توسط دولت یا مراجع صلاحیت دار که در قوانین و مقررات مشخص هستند، تولید و منتشر می‌شود و اطلاعاتی را در مورد وضعیت عمومی کشور برای امور مدیریتی (برنامه ریزی، سیاست گذاری، و تصمیم‌گیری) به دست می‌دهد. در واقع آمار رسمی به اطلاعات عمومی مربوط می‌شود که به نفع جامعه و با بودجه دولتی تولید می‌شود. آمار رسمی برای همگان قابل دسترس خواهد بود.

**گردآوری داده<sup>۱۰</sup>:** تمام فرایندهایی را شامل می‌شود که مرتبط با بدست آوردن داده از جامعه آمارگیری است.

**به موقع بودن<sup>۱۱</sup>:** به فاصله زمانی بین تاریخ پیش‌بینی شده (اسمی) و تاریخ واقعی انتشار نتایج گفته می‌شود.

**بهنگام بودن<sup>۱۲</sup>:** به فاصله زمانی بین پایان دوره مرجع و تاریخ واقعی انتشار نتایج گفته می‌شود.

**لایه انتشار<sup>۱۳</sup>:** برای تولید انتشارات و محصولات آماری به صورت بر خط و برون خط (Online/Offline) و تحلیل و گزارش‌گیری است.

**لایه ذخیره‌سازی<sup>۱۴</sup>:** برای ذخیره‌سازی داده‌های و فراداده‌های معتبر و مرتبط است.

<sup>7</sup> Knowledge Managemen

<sup>8</sup> Data Mining

<sup>9</sup> Official Statistics

<sup>10</sup> Data Collection

<sup>11</sup> Punctuality

<sup>12</sup> Timeliness

<sup>13</sup> Dissemination Layer

لایه تولید<sup>۱۵</sup>: برای جمع‌آوری، اعتبارسنجی، پردازش و مدیریت داده و فراداده‌ها است. انبارداده<sup>۱۶</sup>: مجموعه‌ای از داده‌های موضوع‌گرا، یکپارچه و با تغییرات زمانی است که تصمیمات مدیریتی را حمایت می‌کند.

در واقع تعاریف مختلفی برای انبارداده وجود دارد. دو تعریف مهم و شناخته شده برای انبارداده توسط کیمبال و اینمون ارائه شده است که به آن‌ها در ادامه اشاره می‌شود. رالف کیمبال - کارشناس زبده‌ی انبارداده - تعریف خلاصه‌ای از یک انبارداده ارائه داده است. در این تعریف وی بر روی وظیفه‌ی یک انبارداده متمرکز شده اما بیل اینمون به چگونگی ایجاد یک انبارداده نیز پرداخته است. از نظر اینمون، انبارداده، نسخه‌ای از داده‌های تراکنشی است که به‌طور ویژه برای انجام پرس و جو و تحلیل ساخت‌یافته شده است.

یکی از متداول‌ترین تعاریف انبارداده توسط بیل اینمون - دانشمند آمریکایی رایانه - ارائه شده است که بسیاری از صاحب‌نظران وی را به‌عنوان پدر این دانش به رسمیت می‌شناسند:

یک انبارداده، مجموعه‌ای از داده‌ها است که به‌صورتی «موضوع‌گرا»، «یکپارچه»، «متغیر با زمان» و «از دست نرفتنی» گرد آمده‌اند و در پشتیبانی فرایند تصمیم‌سازی‌های مدیریتی جای می‌گیرد.

موضوع‌گرا<sup>۱۷</sup>: یک انبارداده باید بتواند برای تحلیل یک حوزه‌ی موضوعی خاص مورد استفاده قرار بگیرد. برای مثال، «خدمات»، «متقاضیان»، «پیمان‌کاران»... می‌توانند هر کدام یک موضوع خاص باشند. روش‌های زیادی برای طبقه‌بندی یا موضوعی کردن داده‌ها وجود دارد اما چیزی که در پیاده‌سازی و به‌کارگیری یک انبارداده اهمیت دارد این است که تصمیم‌سازان به‌عنوان کاربران نهایی به دنبال تحلیل چه چیز هستند و منابع اطلاعاتی توسط چه بخش(هایی) تولید می‌شود. برای مثال دادگان استاندارد یک سازمان، ممکن است وظایفی را برای ذخیره‌سازی اطلاعات و ام‌ها، کارت‌های اعتباری، پس‌انداز و

<sup>14</sup> Storage Layer

<sup>15</sup> Generation Layer

<sup>16</sup> Data Warehouse

<sup>17</sup> Subject-Oriented

ودیعہ بر عہدہ داشتہ باشد (عملیات‌گرا) اما یک تحلیل‌گر بخواهد اطلاعاتی در مورد مشتریان، فروشندگان، محصولات و فعالیت آن‌ها داشته باشد. جهت‌بخشی به داده‌ها به‌سوی موضوعی خاص، منجر به انجام طبقه‌بندی‌هایی مفید برای یک تحلیل‌گر سازمانی خواهد شد.

یکپارچه<sup>۱۸</sup>: منابع داده‌ای متعدد درون یک انبارداده به‌صورت یک منبع اطلاعاتی یکپارچه در می‌آیند. برای مثال ممکن است منابع اطلاعاتی A و B حاوی روش‌های مختلفی برای شناسایی یک محصول یا ویژگی باشند. اما در یک انبارداده تنها یک روش برای شناسایی یک محصول یا ویژگی وجود خواهد داشت. به‌عبارتی دیگر خاصیت یکپارچگی یعنی با یک مدل واحد برای ذخیره‌سازی داده‌ها روبه‌رو هستیم. از آنجا که داده‌های عملیاتی یک سازمان در منابع گوناگونی ذخیره می‌شوند مدل داده‌ای هر کدام نیز متفاوت خواهد بود. هدف از پیاده‌سازی یک انبارداده این است که کاربر هنگام اجرای یک پرس و جو با مدل‌های داده‌ای مختلف سر و کار نداشته باشد. یکپارچه‌سازی داده‌ها در یک مکان و در یک مدل داده‌ای یکی از خواص اصلی یک انبارداده یا یکی از اهداف پیاده‌سازی آن است.

متغیر با زمان<sup>۱۹</sup>: داده‌های موجود در یک انبارداده مبتنی بر زمان نگهداری می‌شوند و می‌توان داده‌های مربوط به سه، شش، دوازده ماه پیش یا حتی قدیمی‌تر را از آن به‌دست آورد. این خاصیت در تضاد با سیستمی است که به دلایلی جدیدترین داده‌ها در آن نگهداری می‌شود. برای مثال یک سیستم ثبت معاملات ممکن است شامل آدرس‌های اخیر یک مشتری باشد در صورتی که یک انبارداده می‌تواند حاوی تمامی آدرس‌های مرتبط با آن مشتری باشد. داده‌های موجود در یک انبارداده می‌توانند در برگیرنده‌ی اطلاعات مربوط به یک دهه باشد در حالی که دادگان استاندارد معمولاً فقط شامل داده‌های یک ماه گذشته هستند. در یک انبارداده، زمان یک عامل کلیدی است چون به کاربر اجازه می‌دهد تا برای تحلیل خود یک روند در نظر داشته باشد. برای مثال با چنین

---

<sup>18</sup> Integrated

<sup>19</sup> Time-Variant

قابلیتی می‌توان یک ویژگی را برای یک فصل از سال استخراج و مورد بررسی قرار داد و آن را دوره‌های زمانی دیگر مقایسه کرد.

از دست نرفتنی<sup>۲۰</sup>: هنگامی که داده‌ای در انبار داده قرار می‌گیرد تغییر نخواهد کرد. بنا بر این داده‌های موجود در یک انبار داده که مبتنی بر تاریخ هستند هرگز نباید تغییر یابند. در یک دادگان استاندارد، داده‌ها می‌توانند اضافه، حذف یا روزآمد شوند و معمولاً حالت فعلی سازمان را نشان می‌دهند. در یک انبار داده، داده‌های جدید فقط می‌توانند اضافه شوند و کاربران نمی‌توانند داده‌های موجود در آن را تغییر دهند. در انبار داده برخی سازمان‌ها حتی اطلاعات ۱۰ تا ۲۰ سال ذخیره می‌شود.

انبار داده موضوعی<sup>۲۱</sup>: یک انبار داده موضوعی انبار داده‌ای است که برای یک واحد ویژه در یک سازمان ایجاد شده است. برای مثال واحدهای مالی، فروش و بازاریابی یک سازمان می‌توانند انبار داده موضوعی خود را داشته باشند. هم‌چنین تعریف دیگری برای انبار داده موضوعی به این صورت بیان می‌شود که یک انبار داده موضوعی یک لایه‌ی دسترسی به محیط انبار داده است و از آن برای استخراج داده برای کاربران استفاده می‌شود. به عبارت دیگر انبار داده موضوعی زیرمجموعه‌ای از یک انبار داده است که برای گروه مخصوصی در کسب و کار مورد نظر جهت‌دار شده است. عملاً انبار داده موضوعی قسمت‌های کوچکی از یک انبار داده هستند.

مرکز داده<sup>۲۲</sup>: مرکز داده، مجموعه‌ی بزرگی از سرورهای شبکه‌های رایانه‌ای است که معمولاً برای ذخیره‌سازی، پردازش یا توزیع حجم زیادی از داده‌ها از راه دور توسط سازمان‌ها مورد استفاده قرار می‌گیرد.

تلخیص<sup>۲۳</sup>: تلخیص، روشی برای افزایش سرعت اجرای پرس و جوها است. در این روش اطلاعات مرجع با توجه به ابعاد انتخاب‌شده توسط کاربر خلاصه می‌شوند.

استخراج؛ تبدیل و بارگذاری<sup>۲۴</sup>: فرایندی است که طی آن داده‌ها از یک محیط به محیط دیگر جابجا می‌شوند.

<sup>20</sup> Non-Volatile

<sup>21</sup> Data Mart

<sup>22</sup> Data Center

<sup>23</sup> Aggregation

<sup>24</sup> Extract, transform and load