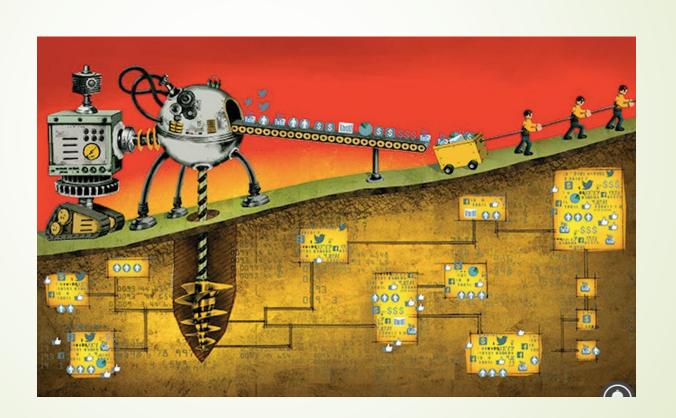
# داده کاوی

# Data mining

آرزو حبیبی راد گروه آمار دانشگاه فردوسی مشهد

# بخش اول:

## معرفی داده کاوی و ارتباط ان با سایر علوم



### مقدمه

- چرا دادهکاوی؟
- مگر در آمار با دادهها سرو کار نداریم؟
- چرا برای تحلیل داده ها، دادهکاوی مطرح می شود؟
- داده کاوی ترجمهی کلمه «Data Mining» بوده که به معنای «کاویدن معادن داده» است.
  - داده کاوی یعنی استخراج اطلاعات گرانبها از حجم عظیم معادن داده!
- کلمه Mining در معنای تحت الفظی خود یعنی «استخراج از معدن» است و در واقع عبارت Data Mining نشان می دهد که حجم انبوه اطلاعات مانند یک معدن عمل می کند و از ظاهر آن مشخص نیست چه عناصر گرانبهایی در

عمق این معدن وجود دارد. تنها با کند و کاو و استخراج این معدن است که می توان به آن عناصر گرانبها دست پیدا کرد.

### دادهکاوی چیست؟

- ◄ داده کاوی یک فرآیند محاسباتی است که در و اقع الگو یا الگو هایی را در مجموعه از داده های عظیم کشف میکند. در تمامی تعریف های مرتبط به داده کاوی کلمه کشف کردن را میتوان پیدا کرد.
  - ✓ داده کاوی شاخه ای از علوم کامپیوتر و ریاضی است که در واقع ترکیبی از تکنیکهای آماری، علوم اطلاعات، یادگیری ماشینی و نظریه پایگاه داده است.
- ✓ در سال ۱۹۸۹ اصطلاح کشف دانش از پایگاه داده (Knowledge Discovery form Databases) یا به اختصار (KDD) مطرح شد و در سال ۱۹۹۰ اصطلاح داده کاوی به وجود آمد و در همین سال با استفاده از نراده کاوی، خرده فروشها و بازارهای مالی به تجزیه تحلیل دادهها و پیش بینی نوسانات در نرخ بهره و افزایش مشتری پرداختند.



• در دهه 1930 اولین بار و اژه داده کاوی به کار گرفته شده است، و شکوفایی آن همزمان با رشد ساختار پایگاه داده ای صورت گرفت.

• داده کاوی در پاسخ به نیاز «تحلیل داده های بزرگتر از حد و اندازه معمول برای کشف دانش و الگو » شکل گرفت جایی که الگوریتمها و روشهای آماری از لحاظ تئوری پاسخگو نبودند، مفاهیم داده کاوی پا به عرصه میگذاشتند.

✓ نکته مهم در دادهکاوی این است که کشف دانش، عموما توسط یک خبره و در واقع یک انسان صورت گرفته و می گیرد.

• حوزهها و محدوده آن شامل همه روشهای داده محور تحلیلی میباشد و زمانی از دادهکاوی استفاده میشود که مسئله ای بروز پیدا کرده و کاربر و خبره قصد حل آن را داشته باشد. بنابراین داده کاوی «مسئله محور» است.

### چرا دادهکاوی؟

- با رشد و افزایش توجهات به دادهکاوی، پرسش «چرا دادهکاوی؟» همواره مطرح می شود در پاسخ به این پرسش باید گفت، دادهکاوی دارای کاربردهای زیادی است. بدین ترتیب، زمینه ای جوان و آینده دار برای نسل کنونی محسوب می شود.
- این زمینه تو انسته توجهات زیادی را به صنایع و جوامع اطلاعاتی جلب کند. با وجود گستر و وسیع داده ها، نیاز حتمی به تبدیل چنین داده هایی به اطلاعات و دانش وجود دارد.
- بنابراین، بشر از اطلاعات و دانش برای گستره وسیعی از کاربردها، از تحلیل بازار گرفته تا تشخیص بیماریها، کشف کلاهبرداری و و پیش بینی قیمت سهام استفاده میکند.
- در مجموع باید گفت، ضرب المثل انگلیسی « نیاز مادر همه ابداعات بشر است»، پاسخی کوتاه و گویا به پرسش مطرح شده است.

# دادهکاوی و علم آمار

- داده کاوی شاخه ی توسعه یافته و پیچیده ی علم آمار است. علم آمار به تنهایی برای صاحبان صنایع و شرکتها بسیار اهمیت دارد.
- تفاوت اصلی علم آمار با شاخه ی توسعه یافته ی خود یعنی داده کاوی، در حجم داده های مورد تحلیل و روش مدلسازی دادههاست. در بیشتر موارد حجم داده های اولیه ی مورد استفاده در داده کاوی آنقدر زیاد است که به یک چالش زمان بر و پر هزینه تبدیل می شود. اما مدلسازی داده های و رودی و دستیابی به اطلاعات پنهان و ارزشمند موجود در این حجم عظیم داده، با کمک هوش مصنوعی و شیوههای خودکار یادگیری انجام می شود.
  - در داده کاوی به دنبال ار تباطهای پنهان و کشف نشده هستیم به طوری که داده کاوی پل ارتباطی میان علم آمار، علم کامپیوتر، هوش مصنوعی، یادگیری ماشین داده می باشد.
- وظیفه ی داده کاوی، کاویدن و استخراج دانش از منابع عظیم داده است تا اطلاعات گرانبهایی که در حجم انبوهی از اطلاعات سطحی پنهان شده است را آشکار سازد.

# ارتباط داده کاوی با یادگیری ماشین، هوش مصنوعی و علم داده

#### الله مصنوعي:

- ✓ هوش مصنوعی منظور دنیایی است که در آن رباتها بر انسانها چیره میشوند و حیات و ممات نسل بشر در معرض تهدید قرار میگیرد.
- ✓ اما در علم کامپیوتر، هوش مصنوعی استفاده از الگوریتمهایی است که بتوانند رفتار انسانی را در قالب کدهای کیفر
  و یک ریاضی در آورده و رفتار مناسب با آن رفتار صورت گرفته را ارائه دهند.
- ✓ حال این شبیه سازی رفتار انسانی می تواند در قالب الگوریتم و منطق ریاضی «اگر آنگاه»، یا در قالب اینکه خود ماشین و ربات یاد بگیرد، باشد.
- ✓ هرچه که هست، هوش مصنوعی به رباتها قابلیت تصمیم گیری میدهد و در واقع اینجا خبره به جای استفاده از الگوریتمها برای حل یک مسئله (که در داده کاوی مورد بحث است)، یک بار برای همیشه کلیه حالتهایی که به ذهناش میرسد را در قالب قواعد و منطق به ماشین تعریف می کند یا کاری می کند که ماشین خودش یاد بگیرد.

### یادگیری ماشین:

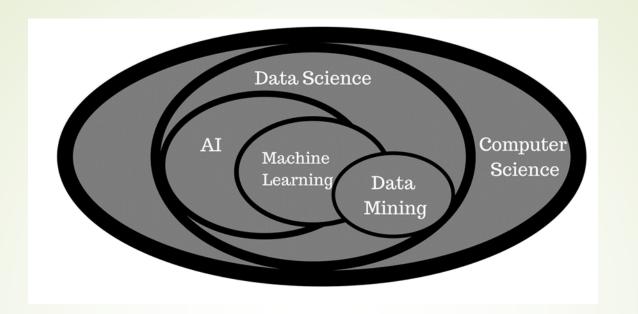
- ✓ یادگیری ماشین بخشی از هوش مصنوعی است.
- ✓ بخشی که ماشین ها به جای تطبیق دهی میان قواعد از پیش تعریف شده، خود یاد می گیرند و برای حالت های جدیدی که مشاهده نشده، تصمیم می گیرند. با این نگاه، یادگیری ماشین بخش اصلی هوش مصنوعی است و هر روز این اهمیت بیشتر می شود چرا که دیگر نمیتوان تمام قواعدِ رفتار های پیچیده انسانی را کشف کرد و به خورد ربات داد بلکه می توان با داده های زیادی که این روز ها تولید می شود، کاری کرد که ربات و ماشین خود یاد بگیرد.
  - ✓ با این بحث ها، می توان گفت یادگیری ماشین و هوش مصنوعی قرین هم بوده و بسیار شبیه به هم هستند.

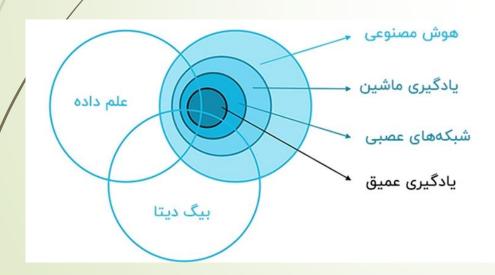
# معرفی علم داده:

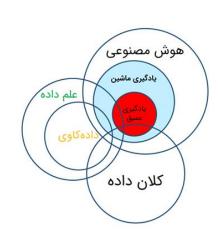
✓ علم از هر نوع آن به بررسی و کشف روابط موجود در آن زمینه می پردازد. برای مثال، علم فیزیک به کشف قواعد دنیای ماده می پردازد. همین طور علم ریاضی به ارائه تعریفهای گوناگون از منظر های گوناگون برای پدیده های طبیعی از نظر عددی (و نه لزوما داده) می پردازد.

✓ علم داده، نیز فارغ از اینکه کاربردها آن در چه جاهایی مصداق پیدا میکند، به دنبال شناخت ماهیت خود داده، الگوریتمهای موجود و ایجاد مفاهیم جدید برای کارایی بیشتر داده در تئوری و عمل است.

✓ لذا هر جا قرار بر این باشد که الگوریتم جدیدی در حوزه داده ایجاد شود، ترجیح به کار بردن واژه «علم داده » است.







# چه مسئلهای علت نیاز به داده کاوی شده است؟

- √ میتوان گفت اصلی ترین دلیل توجه به این دانش در اختیار داشتن حجم وسیعی از داده ها موجود است، که این نیاز به وجود آمده بود که باید از این داده ها اطلاعات و دانش مفیدی استخراج کرد و از آن برای کنترل، تحلیل و پیشبینی استفاده کرد.
  - ✓ بنابراین دلایل کلی را میتوان در سه مورد خلاصه کرد:
    - 1. حجم داده ها با سرعت زیادی در حال رشد است.
      - 2. اطلاعات ما در مورد این داده ها کم است.
      - 3 دانش ما نسبت به این اطلاعات صفر است.

- ✓ اما پاسخ تمامی این مشکلات تنها و تنها داده کاوی است.
- ✓ با دادهکاوی میتوان حجم عظیم داده ها را اصطلاحاً کاوید و دانش را از دل آن بیرون کشید.

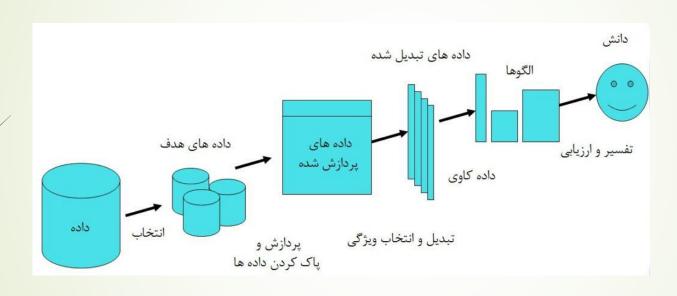
#### ن ویژگیهای اصلی دادهکاوی:

- کشف اتوماتیک الگوها در انجام پروژههای داده کاوی پیش بینی احتمالی نتایج و خروجی ها
  - تولید اطلاعات اجرایی و مفید متمرکز بر دادههای بزرگ.

#### ♦ مزایای دادهکاوی:

- ایجاد روابط بصورت اتوماتیک
  - استفاده از دادههای متنوع
    - دینامیک بودن
  - عدم نیاز به دادههای صحیح
    - ساخت مدلهای واقعی
  - تحلیل کردن دادههای و اقعی
- دوری از اشکالات احتمالی نمونهگیری

### مراحل دادهکاوی:



- 1. پالایش دادهها
- 2. يكپارچەسازى دادەھا
  - 3 انتخاب دادهها
    - 4 تبدیل دادهها
      - 5. دادهکاوی
- 6. ارزيابي و تفسير الگوها
  - 7. ارائه دانش

# چهار گام اصلی دادهکاوی:

### > گام صفر: تعریف مسئله شناخت دادهها

قبل از شروع فرآیند داده کاوی به تعریف دقیق مساله اختصاص پرداخته می شود و روشهای به کارگیری دادهها جهت پاسخگویی به مسئله مدنظر قرار می گیرد. در این گام باید به خوبی داده ها را شناخت.

#### → گام اول: آمادهسازی دادهها (پالایش دادهها)

در این گام، تحلیلگر ممکن است رو شهای مختلف دادهکاوی را بر روی دادههای آماده شده امتحان کند تا بتواند به مدف اصلی پروژه برسد.

ساخت مدل یک فرآیند خطی نیست و رفت وبرگشتهای زیادی وجود دارد. یک مدل بهینه در داده کاوی وجود ندارد و بسته به مسئلهای که تحلیلگر با آن مواجه است، روشهای مختلف باید آزمایش شوند و خروجی آنها باهم مقایسه گردند. در این مرحله احتمالاً لازم است به گام قبلی بازگشت و برای برخی از الگوریتمها دادهها را به شکل دیگری آماده کرد.

#### کام سوم: یادگیری مدل

بسته به نیاز کسب و کار ، داده کاوی ممکن است با هدف پیشبینی (Prediction) پیدا کردن روابط (Association) و یا برای خوشه بندی (Clustering) استفاده گردد. در هریک از این دسته ها الگوریتم های متفاوتی و جود دارند که بسته به شرایط یکی از آن ها یا ترکیبی از آنان استفاده می شوند.

#### > گام سوم: ارزیابی مدل

در این فرآیند، الگوهای حاصل شده در گام قبل، از جنبههای گوناگونی شامل دقت، صحت و قابلیت تعمیم و دیگر موارد مورد ارزیابی قرار میگیرند. مدلهای ساخته شده تست و ارزیابی میشوند و بهترین مدل از نظر مسئله ی طرح شده در این مرحله، انتخاب میشود. سپس در تبادل نظر با کارفرما، موثر بودن مدل انتخاب شده بررسی میشود. در صورتی که مدل انتخاب شده کمکی در حل مسئله نمی کند کل فرآیند از مرحله ی یک دوباره انجام میشود.

### ح گام چهارم: تفسیر مدل

✓ تفسیر نتیجه در این مرحله نتایج و الگوهای ارائه شده توسط ابزار داده کاوی مورد بررسی قرار گرفته و نتایج مفید معین میشود.

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

### در دادهکاوی چه دادههایی کاوش میشوند؟

ح در داده کاوی داده های ثبت شده در پایگاه داده ها یا انبار داده ها مورد تحلیل و تفسیر قرار می گیرند.

#### پایگاه دادهها:

- ✓ مجموعه ای از داده های ثبت شده توسط سیستمهای کاربردی در مؤسسات، سازمانها، مراکز تجاری، صنعتی، ... را پایگاه داده می گویند.
- ✓ دادههای یک سیستم پایگاه داده شامل مجموعهای از دادههای مرتبط و نرمافزارهایی جهت دستیابی به این دادهها میباشد.
  - ✓ پایگاه داده ها توسط متخصصین علم کامپیوتر نوشته و آماده میشوند.
  - ✓ مثالهایی از پایگاه داده: ثبت اطلاعات افراد توسط کارمندان یک فروشگاه زنجیرهای.

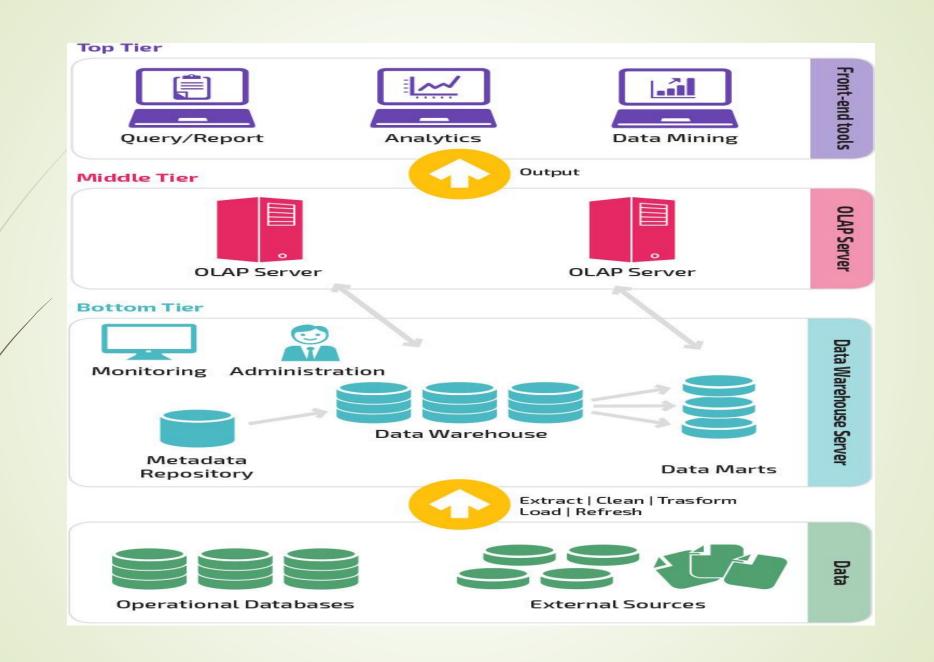
### انبار دادهها:

- ✓ مخزنی از داده ها هستند که از چندین پایگاه داده جمع آوری شده اند و به صورت یک شمای یک دست و کلیاند و غالباً در سایت ذخیره میشوند، را ( Data warehousing) گفته نمی شود.
  - ✓ با کمک انبار داده ها می توان دانش یا اطلاعات جامع تری کسب نمود.
- ✓ در این انبار، اطلاعاتی ذخیره می شود که یک ساز مان برای تصمیم گیری های را هبردی و کلان به آن نیاز دارد.
  - ✓ به فرایند ساخت و استفاده از انبار داده ها انبارش داده ها می گویند.
  - ✓ پایگاه داده ها با کمک انبارش داده تبدیل به یک منبع مفید برای تحلیل داده ها می شود.
- ✓ امروزه موفق ترین شرکتها آنهایی هستند که بتوانند به سرعت و با انعطاف بالا به تغییرات و فرصتها پاسخ دهند. کلید پاسخ به این مسئله کاربرد مؤثر و کار آمد داده ها و اطلاعات از سوی تحلیلگران و مدیران است.
   یک انبار داده مخزنی از داده های تاریخی است که بر اساس موضوع سازمان یافته اند تا به تصمیم گیران سازمان کمک کند. زمانی که داده ها در یک دادهگاه یا انبار داده ذخیره می شوند، می توان آن ها را مورد ارزیابی قرار داد.

- به عبارتی انبار داده ها مجموعه ای از داده های موضوع –گرا، یکپارچه شده، زمان-متغیر و غیر فرّار هستند.
- موضوع-محور: دادهها در انبار داده بر خلاف سیستم های عملیاتی پیرامون موضوعات سازمان گرد
   هم می آیند. جهت گیری موضوعی برای تصمیم گیری واقعا مفید است. گردآوری اشیای مورد نیاز،
   موضوع-محوری نامیده می شود.
  - یکپارچه شده: دادههای موجود درون یک انبار داده یکپارچه هستند. از آنجا که این دادهها از چند سیستم عملیاتی میآیند همه ناسازگاریها میبایست حذف شوند.
- خصوصیات آن شامل قراردادهای نام گذاری، اندازهگیری متغیرها، ساختار داده، خصوصیات فیزیکی داده و ُ مواردی از این دست است.
- زمان-متغیر: با این که سیستمهای عملیاتی به دلیل پشتیبانی از عملیاتهای هر روزه مقادیر فعلی را نشان میدهند؛ اما انبارهای داده، نشان دهنده دادههایی در افق زمانی بلندمدت وستند. این بدان معنی است که انبار داده حاوی دادههای تاریخی است. انبار داده به طور عمده برای دادهکاوی و پیشبینی استفاده میشود، چون اگر کاربری در جستجوی یک الگوی خرید برای یک مشتری خاص باشد، میبایست به دادههایی در مورد خریدهای فعلی و گذشته نگاه کند.
  - **عیر فرّار:** دادههای موجود در انبار داده تنها خواندنی هستند، یعنی نمیتوان آنها را به روزرسانی، ایجاد یا حذف کرد.

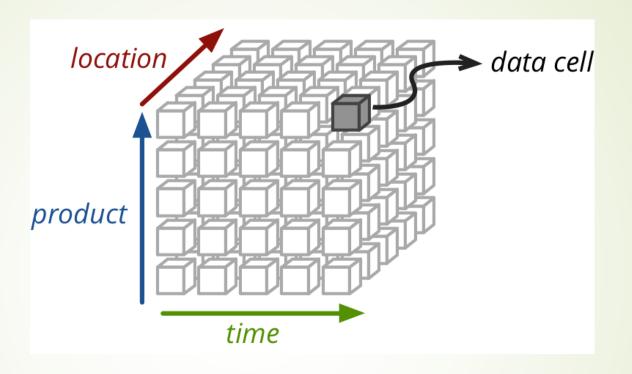
- √ در انبارش داده ها با دو سیستم OLAP و OLTP سروکار داریم.
- On Line Transition Procedure مخفف عبارت OLTP
- On Line Analytically Procedure مخفف عبارت OLAP واژه
- ✓ سیستم OLTP یک سیستم مشتریگرا است، کاربران این سیستم معمولا کارمندان سازمان هستند.
- ✓ سیستم OLAP یک سیستم بازارگرا است، کاربران این سیستم معمولا متخصصین و مدیران سطح بالای سازمان هستند.
  - ✓ طراحی انبارش داده و لایه های مختلف معماری انبارش داده:

انبارش داده ها ساخت یک معماری چند ر دیفی بر روی داده ها می باشد. غالبا از معماری سه ر دیفی استفاده می شود.

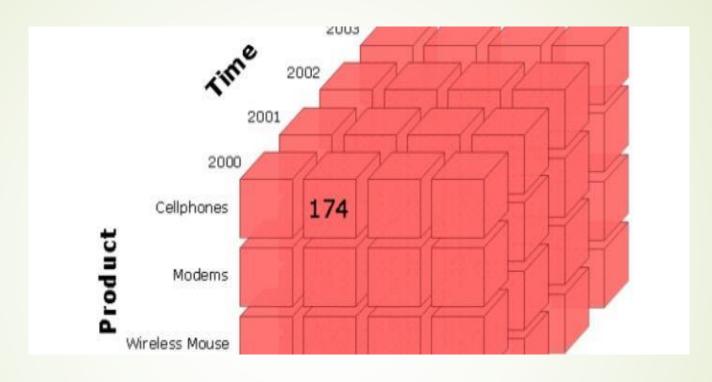


# مكعب دادهها (Data Cube)

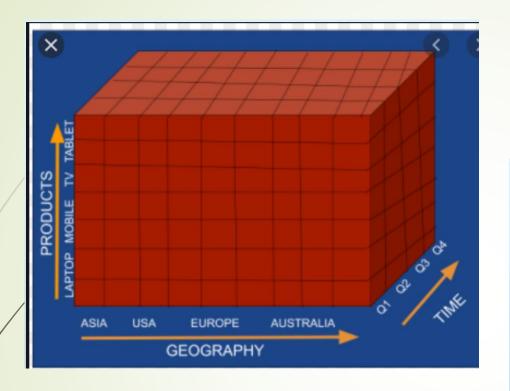
- > انبار داده ها بر اساس یک مدل داده ای چند بعدی به نام مکعب داده نمایش داده می شود.
- ح مکعب داده، نمایش یا الگو چند بعدی از داده ها است به طوری که این امکان را فراهم میکند تا داده ها در چند بعد نمایش داده شوند.
- مکعب داده چند بعدی (MOLAP): اکثر محصولات OLAP بر اساس ساختاری طراحی شدهاند که داده ها را به شکل آرایه های مرتب شده بر اساس ابعاد داده ذخیره میکنند.
  - ✓ MOLAP در مقایسه با رویکردهای دیگر و به دلیل ساختار خاص نگهداری اطلاعات، از سرعت بالایی برخوردار است و کارایی را بهبود میبخشد. وقتی که تعداد ابعاد بزرگتر است، مکعب کوچکتر می شود.
  - مکعب داده رابطه ای (ROLAP): OLAP ارتباطی یا ROLAP، از مدل پایگاه داده رابطهای استفاده میکند و عملیات OLAP را در قالب استفاده از مدل Relational ساماندهی میکند و لذا از سیستمهای MOLAP کندتر عمل میکنند و به فضای بیشتری نیز نیاز دارند.

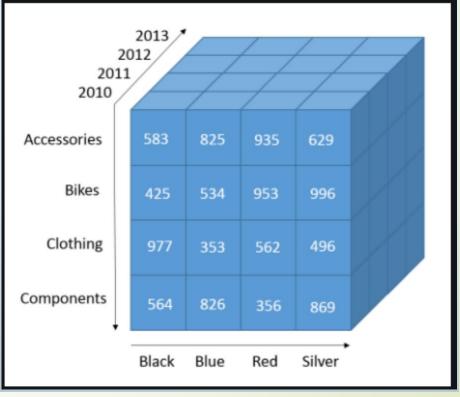


❖ یک انبار داده معمولا با کمک یک ساختار چند بُعدی با نام مکعب داده مدل سازی میشود. به عبارتی یک مکعب داده، یک دید چند بُعدی را به کاربر ارائه میکند و با پیش محاسبه ی آن می توان دسترسی سریعی به داده های خلاصه شده داشت.

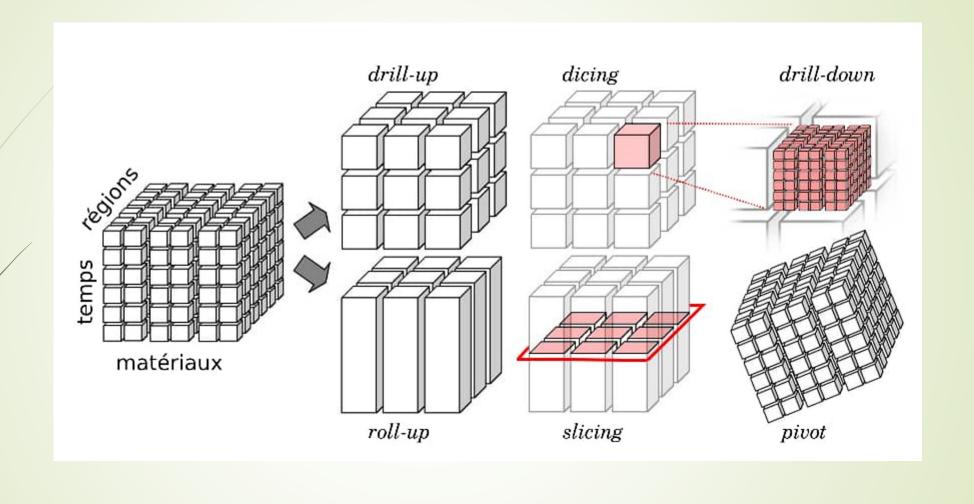


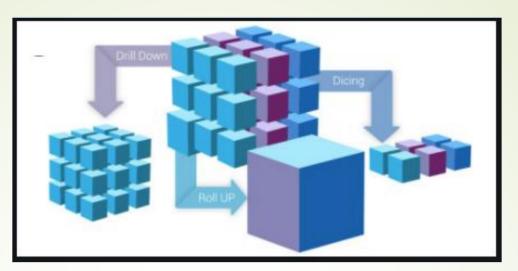
یک انبار داده معمولا با کمک یک ساختار چند بُعدی با نام مکعب داده مدل سازی
می شود. به عبارتی یک مکعب داده، یک دید چند بُعدی را به کاربر ارائه می کند و
با پیش محاسبه ی آن می توان دسترسی سریعی به داده های خلاصه شده
داشت.

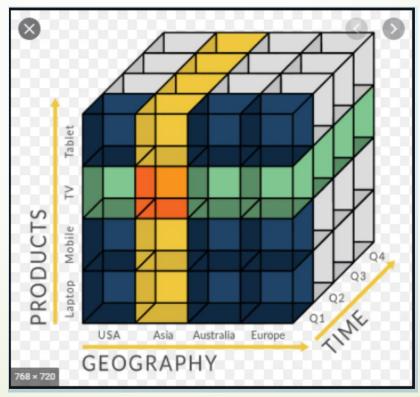


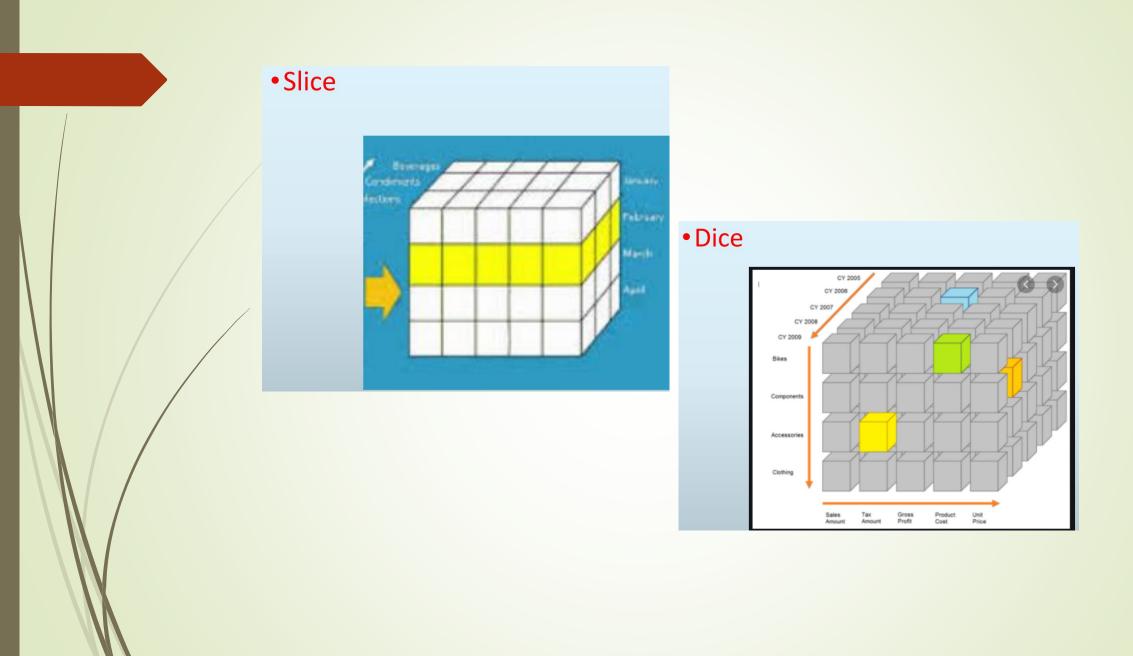


- پنج عمل اصلی زیر را میتوان بر روی مکعب داده ها انجام داد:
- 1. ROLL-UP: با کاهش تعداد بعدها ساختار جدیدی را در مکعبها ایجاد میکند.
- 2. DRILL-DOWN: برعکس عمل ROLL-UP است و اطلاعات را به صورت مفصلتری نشان میدهد. برای مثال زمانی که به جای اطلاعات فروش فصلی، اطلاعات فروش ماهانه مورد نیاز باشد از این عمل استفاده میشود.
  - 3. SLICE: انتخابی را بر روی یک بعد انجام میدهد.
- 4. DICE: انتخاب را روی چند بعد انجام میدهد به عبارتی مکعبها را به صورت جدا جدا بیرون می آورد.
  - 5. PIVOT: با چرخش مکعب از یک بعد دیگر میتوان به آن نگاه کرد (جای محور ها را میتوان عوض کرد).









بخش دوم:

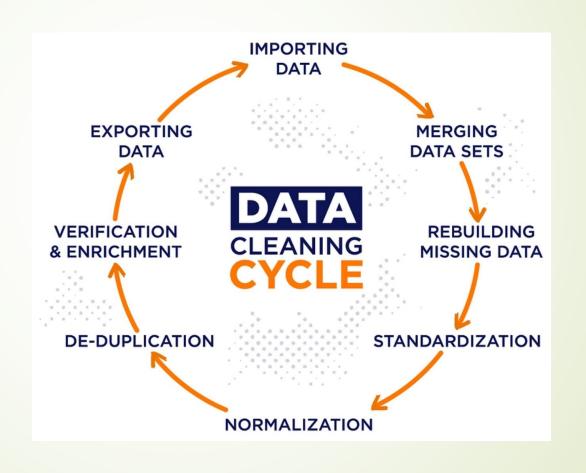
آماده سازی و پیش پردازش داده ها



در این مرحله هدف آماده کردن داده ها برای قسمت تحلیلی با استفاده از روش های داده کاوی است و معمولا بیشترین زمانی که در انجام پروژه داده کاوی صرف می شود در این مرحله انجام می شود. به عبارتی %۸۰ زمان یک پروژه مرتبط به آماده سازی داده است و دلیلش هم آن است که در دنیای واقعی داده ها آنطوری نیستند که باید باشند! وجود نمادهای نامربوط، داده های پرت، گمشده، ناسازگاری و ... دلیلی برای زمان زیادی است که در این مرحله هدر می رود و زمان هایی نیز پیش می آید که داده ها به صورت الکترونیکی ذخیره نشدند و برگردان این اطلاعات به صورت الکترونیکی ذخیره نشدند و برگردان این اطلاعات به صورت الکترونیکی کاری بسیار طاقت فرسا و زمان بر است.



آماده سازی داده ها وقت زیادی را در داده کاوی به خود اختصاص میدهد. آماده سازی داده ها ضروری می باشد.



### ✓ آمادهسازی دادهها شامل سه مرحله به شرح زیر است:

- 1. شناخت انواع دادهها
- 2. پیشپردازش دادهها
  - 3. پویش دادهها
- این مرحله از دادهکاوی به تجمیع و پاکسازی دادههایی اختصاص مییابد که در مرحله تعریف مسئله شناسایی شدهاند.
  - ممکن است داده ها در سر اسر ساز مان توزیع شده و در قالب های مختلفی ذخیره شوند و یا ممکن است شامل تناقضات و ناسازگاری هایی از جمله و رودی های نادر ست یا گمشده باشد.
- فرآیند پاکسازی داده ها تنها به حذف داده های نامناسب و یا وارد کردن مقادیر از دست رفته خلاصه نمی شود بلکه کشف روابط پنهان شده میان داده ها، شناسایی دقیق تر منابع داده و تعیین مناسب ترین متغیر ها برای استفاده در مرحله تحلیل را در برمی گیرد.

#### دادهها:

- ✓ هر چیزی که قابل ثبت شدن باشد داده است.
- ✓ داده می تواند عددی باشد یا غیر عدد باشد.
- ✓ به هر رقم، عدد، کاراکتر، رشته، سیگنال، ... و به طور کلی به هرآنچه که قابل ثبت شدن باشد داده می گویند.

#### ویژگی دادهها:

- ✓ داده های ثبت شده باید با یک مقیاس مناسب اندازه گیری شودند.
- ✓ منظور از اندازه گیری نسبت دادن یک عدد حقیقی به یک ویژگی میباشد که این امر طبق یک قائده ی مشخص صورت گیرد.

### مقیاسهای اندازهگیری دادهها:

- 1. مقیاس اسمی: تنها برای تمیز میان اسم افراد به کار میرود و نمی توان آنها را برای مقایسه 4 عمل اصلی به کار برد.
- 2. مقیاس ترتیبی: داده ها خاصیت ترتیبی دارند اما قابلیت انجام عملیات حسابی وجود ندارد.
- 3. مقیاس فاصله ای: داده ها خاصیت ترتیبی دارند و قابلیت انجام عملیات حسابی به طور محدود و جود دارد. در مقیاس فاصله ای صفر، صفر حقیقی نیست.
  - **4.** مقیاس نسبتی: داده ها خاصیت ترتیبی دارند و قابلیت انجام عملیات حسابی به طور کامل و جود دارد.

## ♦ متغیرها و ویژگیهای آنها:

هر رکورد ثبت شده از داده ها، شامل چند ستون داده است که به این ستون ها فیلد یا متغیر می گویند.

### ویژگی متغیرها:

- متغیر کمی یا کیفی
- متغیر گسسته یا پیوسته
- متغیر مستغل یا وابسته

# پیشپردازش دادهها:

❖ داده ها زمانی که نیاز های خواسته شده را برآورده سازند دارای کیفیت هستند. عوامل متعددی و جود دارند که با کمک آن ها کیفیت داده ها توصیف میشوند. مانند صحت و درستی، کامل بودن، بدون نویز بودن، قابل تفسیر بودن و ...

در دنیای واقعی با توجه به حجم زیاد داده و اینکه داده ها از منابع مختلفی جمع آوری می شوند. این داده ها مستعد نویز، ناسازگاری و ناموجود بودن هستند. برای بهبود نتایج کاوش در قسمت اماده سازی داده ها به پیشپر دازش داده ها می پر دازیم.

### روشهای پیشپردازش دادهها در 4 گروه زیر سازماندهی میشوند:

- 1. پالایش دادهها
- 2. يكپارچەسازى دادەھا
  - 3. كاهش دادهها
  - 4. تبدیل دادهها

# يالايش دادهها:

❖ قبل از شروع مدلسازی و تحلیل داده ها انجام پالایش داده ها ضروری است تا داده های مشکل دار، غیر مفید، ....
 مشخص و در مورد آن ها تصمیم گیری صورت گیرد.

مشکلات موجود در دادهها و راههای برخورد با آنها:

■ دادههای گمشده: داده هایی که به دلایل عمدی یا غیر عمدی ثبت نشده اند. در داده کاوی به دلیل حجم زیاد داده، احتمال برخورد با داده های گمشده زیاد است.

- > در صورت مواجهه با دادههای گمشده میتوان از روشهای زیر استفاده کرد:
  - ✓ حذف داده های گمشده
  - ✓ جایگزینی مقدار مناسب
  - ✓ پیش بینی برای خانه های خالی
    - دادههای پرت:

دادههایی که نسبت به اکثر دادهها دچار انحراف یا تمایز زیادی هستند یا به عبارت دیگر خارج از محدوده منطقی دادهها قرار گرفتند.

- > راههای تشخیص دادههای پرت:
  - ✓ روش خوشه بندی داده ها
- ✓ رسم نمودارهای پراکنش یا هیستوگرام

#### روش برخورد با دادههای پرت:

- ✓ ممکن است با مشورت مدیر پروژه حضور داده های پرت در فایل داده ها لازم باشد.
  - ✓ در غیر این صورت مشابه دادههای گمشده برخورد میکنیم.

### دادههای ناسازگار:

- ✓ هرگونه تخریب در مقدار داده که به صورت عمدی یا غیر عمدی صورت پذیرد و منجر به این شود که ما به اصل داده دسترسی نداشته باشیم نویز گفته می شود.
  - ✓ داده هایی که در صورت عمل نویز اتفاق بیافتند داده های نویز دار یا ناسازگار می گویند.
    - ✓ شناخت این نوع داده ها بسیار مشکل و زمان بر هست.
    - ✓ این مشکل در هنگام یکپارچه سازی پایگاه داده ها خیلی رخ می دهد.

### > روش برخورد با دادههای ناسازگار:

## الف) روش هموار سازی یا بسته بندی (Binning):

در این روش پس از تشخیص متغیر نویز دار، ابتدا مشاهدات مرتب میشوند سپس به درون بسته های k تایی قرار می گیرند هر بسته با کمک مقداری براساس میانگین بسته، میانه بسته، کرانه های بسته، ... هموار میشود.

### ♦ روش رگرسیون:

در صورت برخورد با هر نوع اغتشاش در داده ها استفاده از مدل های رگرسیونی برای جایگزینی بهترین راه است. در این روش متغیر هایی که پیشبینی کننده های خوبی برای متغیر و ابسته هستند را انتخاب کنیم و معادله رگرسیون را بدست آوریم. سپس با تعیین مقادیر اولیه برای نقاط پرت با استفاده از خط رگرسیون دوباره رگرسیون را تکرار میکنیم تا مرحله ای که تعداد نقاط مخشوش پیدا شود و با نظر فرد خبره تعیین شود.

# یکپارچهسازی دادهها:

□ در پروژههای داده کاوی ممکن است داده ها از منابع مختلف گردآوری شوند لذا لازم است قبل از تحلیل پایگاههای داده یکپارچه شوند که برای این منظور باید عملیات یکسانسازی ساختار داده ها صورت گیرد.

#### نکات لازم برای یکپارچهسازی دادهها

- ✓ شناسایی موجودیت داده ها (متا داده): شامل جو اب سو الاتی درباره داده است که باعث سهولت در فهم متغیر ها و چگونگی استفاده از داده میشود.
  - نسبت متا داده به داده را میتوان مانند نسبت شناسنامه به فرد در نظر گرفت.
    - ✓ افزونگی و تحلیل همبستگی:
- افزونگی: اگر دو ستون در پایگاه داده ها معرف یک موجودیت باشد با افزونگی مواجه هستیم. ممکن است یک متغیر مانند در آمد سالانه افزونه تلقی شود و اگر بتوانیم با کمک متغیر های دیگر آن را بدست آوریم. برخی افزینگی همبستگی تشخیص داد.

✓ متغیرهای تکراری: رکوردهایی هستند که بار اطلاعاتی جدیدی ندارند و اصلاعات تکراری زیادی در آنها وجود دارد.

• در دادهکاوی به دنبال این هستیم که رکوردهای تکراری را حذف نموده و رکوردهایی ایجاد کنیم که بار اطلاعاتی زیادی دارند. اثر مثبت این کار این است که هم حجم داده ها کم می شود و هم اطلاعات اضافی حذف می شود.

✓ تشخیص تضاد در اندازه گیری مقادیر داده ها: ممکن است مقادیر متغیر مربوط به یک موجودیت در دنیای واقعی در منابع مختلف به صورت متفاوتی نگهداری شده باشد. این تفوات میتواند ناشی از نمایش مقدار یا مقیاسگذاری آن ها باشد.

## کاهش دادهها:

□ تحلیل بر روی دادههای حجیم زمانبر است. به این منظور با افزایش سرعت در تحلیل دادهها میتوان با کمک تکنیکهایی دادهها را کاهش داد که این کاهش میتواند در بعد دادهها، کاهش رکوردها و یا فشر دهسازی باشد.

- > روشهای کاهش دادهها:
  - روش های نمونه گیری
- رسم نمودار های هیستوگرام
- حذف متغییر های غیر ضروری
  - روش های خوشه بندی
    - تحلیل عاملی
    - تجميع داده ها

## تبدیل دادهها:

□ داده ها به گونه ای تبدیل می شوند که فر آیند کاوش ممکن است کار آمدتر اجر ا شوند و الگوهای پیدا شده به راحتی قابل فهم باشد.

- > روشهای تبدیل دادها:
- ✓ ساخت یک متغیر جدید
- ✓ گسسته سازی داده ها ی پیوسته
  - ✓ نرمال سازی داده ها
    - √ هموارسازی
    - ✓ تجميع دادهها

## ساخت یک متغیر جدید:

- در این روش با ترکیب متغیر ها، متغیر جدیدی میسازیم که برای کاوش مورد نیاز است.
  - روش نرمالسازی دادهها:

با نرمال داده ها می توان به گونه ای مناسب داده ها را ر بازه های کوچکتری تبدیل نمود که از جمله آن ها می توان به سه روش زیر اشاره کرد.

#### :min-max روش نرمالسازی

خرض کنید در یک نمونه مینیمم مشاهدات را با  $\chi_{(1)}$  و ماکسیمم مشاهدات را با  $\chi_{(n)}$  نمایش دهیم. برای آنکه هر مشاهده  $\chi_i$  به بازه  $\chi_i$  تبدیل کنیم، لازم است تبدیل زیر را انجام دهیم.

$$x_i' = \frac{x_i - x_{(1)}}{x_{(n)} - x_{(1)}} (b - a) + a$$

### :Z-score روش نرمالسازی

اگر $\overline{X}$  میانگین مشاهدات و S انحراف معیار مشاهدات باشد با تبدیل زیر می توان مقدار مشاهده را به مقدار استاندار دشده تبدیل کرد.

$$Z_i = \frac{x_i - \bar{x}}{s}$$

✓ در صورت وجود داده پرت بهتر است از این روش استفاده شود.

### > روش مقیاسگذاری دهدهی:

در این روش بزرگترین داده به لحاظ قدر مطلق را در نظر میگیریم y)، سپس j را که کوچکترین عددی است به طوری که y را پیدا کرده و تمام اعداد را بر y تقسیم میکنیم که با این کار داده ها در بازه y قرار میگیرند.

#### کسسته سازی دادههای پیوسته:

هدف در گسسته سازی داده ها این است که داده ها با مقیاس های فاصله ای و نسبتی را به داده هایی با مقیاس های اسمی یا ترتیبی تبدیل کنیم.

#### 🔪 تجمیع داده ها:

گاهی اوقات با تجمیع متغیر ها و ساخت متغیر جدید می توان به دادههایی با اطلاعات بیشتر دسترسی داشت، که این امر به آماده سازی

داده ها کمک می کند.

مثال: ساخت متغیر توده بدنی (BMI) با کمک متغیر های وزن و قد

$$BMI = \frac{\mathsf{e}(\dot{\mathsf{c}})}{(\mathsf{e})^2}$$

# پویش دادهها:

هدف اصلی در پویش داده ها این است که پیش از داده کاوی با استفاده از روشهای پویش داده ها، اطلاعات ارزشمندی را از داده ها استخراج کنیم.

### روشهای پویش دادهها عبارتنداز:

- 1) شاخصهای آماری: شامل شاخصهای مرکزی و شاخصهای پراکندگی است
  - 2) مصورسازی دادهها
- 3) پردازش تحلیلی برخط (OLAP): پیشتر در قسمت معرفی مکعب داده به آن اشاره شد

## مصورسازی دادهها:

- ✓ مصور سازی اطلاعات یعنی ارائه داده ها در غالب تصویری و گرافیکی است. این نوع ارائه تصمیم گیران و صاحبین اطلاعات را قادر میسازد، تا آنالیز دادههای خود را به صورت بصری دریافت کنند. همین موضوع باعث می شود تا درک مفاهیم دشوار و شناسایی الگوهای جدید برایشان ساده تر گردد.
  - ✓ ثابت شده که ما با تماشای تصاویر اطلاعات را بهتر درک میکنیم. ما با تجسم همواره یک گام به جلوحرکت خواهیم کرد. نمودارها همواره سریعتر از جداول هستند و نمودارهای چند بعدی که حاصل یک هوش تجاری پیشرفته است می تواند چندین صفحه جدول را در آن واحد به ما نشان دهند.

گرافیک به کمک دانش پایگاه داده ها آمده و حجم عظیمی از اطلاعات را می تواند به ما در یک صفحه، نمایش دهد.

◄ با مصور سازی اطلاعات می توان:

■ 1 - مواضعی که نیاز به توجه و بهبود دارند را برای شما شناسایی کند.

■ 2 - عواملی که رفتار مشتری را تحت تاثیر قرار میدهد، را برایتان نمایان سازد.

■ 3 - به شما کمک کند تا بفهمید چه محصولی در کجا بهتر فروخته می شود.

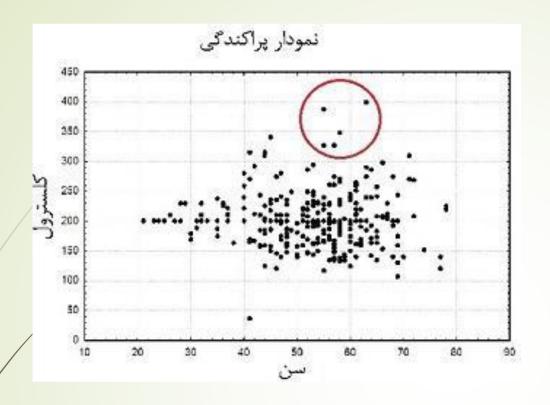
- 4 - پیش بینی (حجم فروش، نقاط فروش و ....) کند.

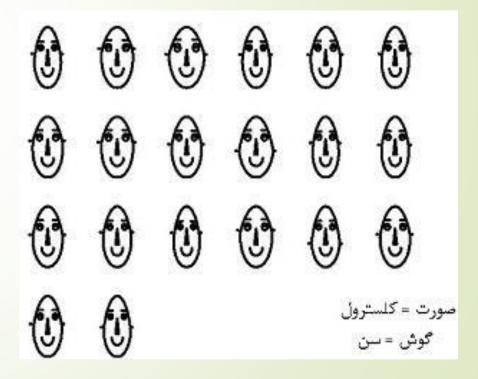


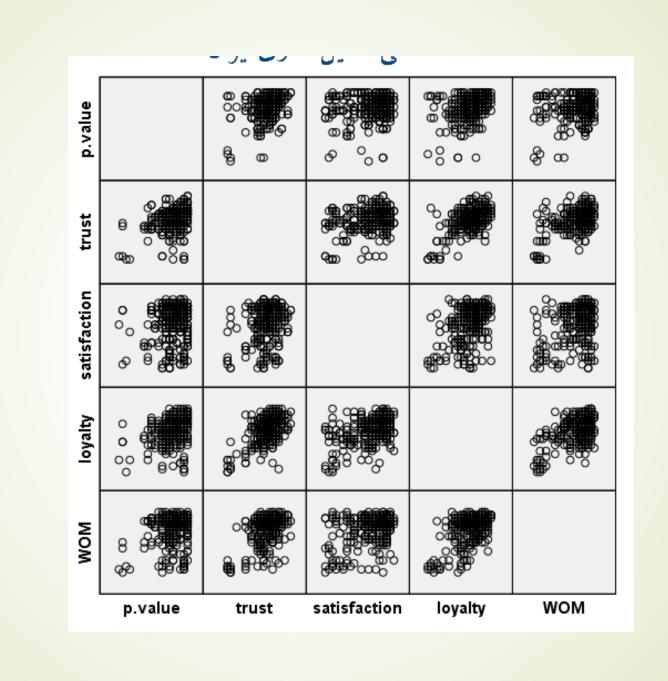
> از روشها مصورسازی دادهها میتوان به موارد زیر اشاره کرد:

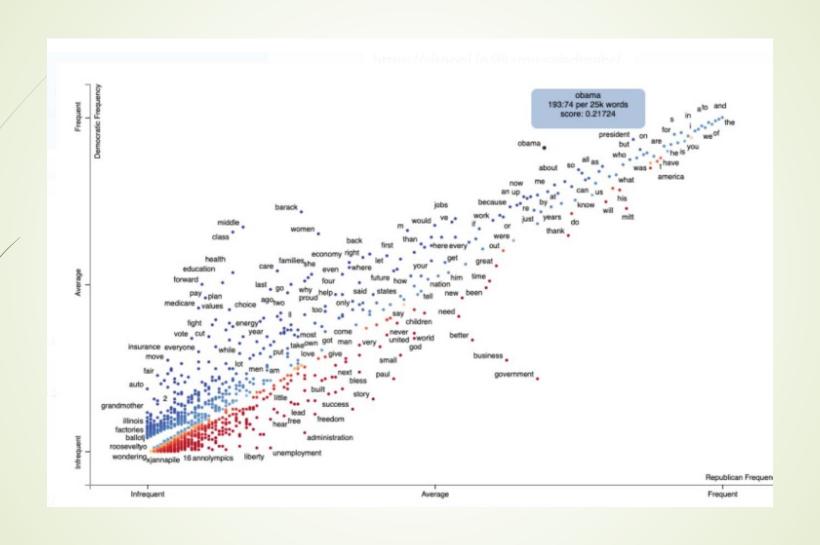
- رسم نمودارها
- مصور سازی پیکسل گرا: در این نوع مصور سازی هر مشاهده به صورت یک پیکسل در صفحه نمایش داده می شود.
- ماتریس نمودار پراکنش: نمودار پراکنش شامل نقاط مربوط به مشاهدات x و Y در مقابل هم می باشد.
  - نمودار وب یا تار عنکبوتی





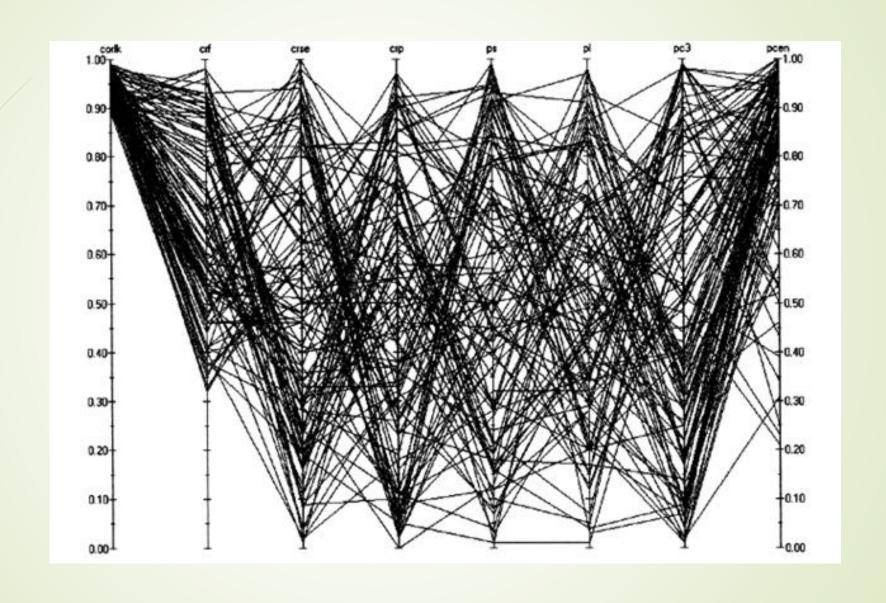






### ✓ نمودار وب یا تار عنکبوتی:

◄ با این نمودار می توان میزان یا به عبارتی شدت روابط بین سطوح مختلف متغیرها را نشان داد. در این نمودار هرچه خطوط بین سطوح پررنگ تر باشد نشان دهنده حجم بیشتر داده ها و یا به عبارتی ارتباط قویتر داده ها در این قسمت بوده و خطوط ظریف تر نشان دهنده ارتباط ضعیف تر می باشد.



## متن کاوی (Text Mining)





