

---

---

### بخش سوم

---

## مدیریت داده کاوی

فصل ششم: انباره داده‌ها

فصل هفتم: متداول‌تری اجرا و پیاده‌سازی پروژه‌های داده کاوی



## فصل ششم

# انباره داده‌ها

سازمانها درک کرده‌اند که سیستم‌های انباره‌داده ابزارهای ارزشمندی در رقابت‌های امروزه هستند. بنگاههای بسیاری، میلیونها دلار برای ساخت انباره‌داده صرف کرده‌اند. طبق تعریف این‌مون<sup>۴</sup> یکی از پیشتازان معماری در ساخت سیستمهای انباره‌داده، انباره‌داده مجموعه‌ای موضوع‌گرا، یکپارچه، از زمانهای مختلف و غیرفرار به منظور پشتیبانی از فرآیند تصمیم‌سازی است. داده‌کاوی چیزی فراتر از پردازش بر روی یک پایگاه داده معمولی می‌باشد. مثالهای زیر این تفاوت را نشان می‌دهند. یک پرس و جوی ساده و پیدا کردن تمامی افراد با نام «علی» در یک پایگاه داده بسیار ساده است ولی در مقابل پیدا کردن افرادی که کارت اعتباری آنها وضعیت مناسبی ندارد و در مرز ورشکستگی می‌باشد، خیلی ساده نیست. پیدا کردن افرادی که بیش از ۱۰۰/۰۰۰ تومان خرید داشته‌اند ساده است ولی در مقابل پیدا کردن افرادی که عادتهای خرید مشابهی دارند و یکسری اقلام خاصی را با هم خرید می‌کنند، کار ساده‌ای نیست. پیدا کردن افرادی که در یک تاریخ خاص از یک فروشگاه خاص شیر خریده‌اند، بسیار ساده است ولی در مقابل پیدا کردن افرادی که غالباً شیر خریداری می‌کنند خیلی ساده نیست.

با توجه به مثالهای مطرح شده، کاملاً مشخص است که داده‌هایی که در داده‌کاوی و خوریتمد.س. ن استفاده بوند، تفاوت عمده‌ای با داده‌های عادی در پایگاه داده‌ها دارند. جهت فراهم کردن این نوع داده‌ها که در انباره‌داده‌ها قرار می‌گیرند، باید یکسری پردازش‌های خاص روی آنها صورت پذیرد. این نوع پردازش‌ها به نام‌های پاکسازی داده‌ها و یکپارچه‌سازی داده‌ها معروفند، که در فصل دوم کاملاً توضیح داده شده‌اند.

## ۶-۱-داده‌کاوی و انباره‌داده‌ها

در دهه ۹۰ میلادی پدیده انباره‌داده‌ها ظهر ریافت. قبل از انبارسازی داده‌ها سیستم‌های کامپیوتری جهت ذخیره، جمع‌آوری، تغییر و تصحیح داده‌ها طراحی شده بودند. این سیستم‌های اولیه به سیستم‌های عملیاتی یا میراثی<sup>۱</sup> موسوم هستند. گرچه جمع‌آوری و ذخیره داده‌ها کارهای مفیدی به حساب می‌آیند، اما دسترسی به آنها و تحلیل داده‌های عملیاتی به راحتی امکان پذیر نیست و علت نیز عدم یکپارچگی داده‌ها می‌باشد. هر برنامه کاربردی داده‌ها را بنا بر نیاز خود تفسیر می‌کند.

مدیران برای تصمیم‌گیری نیازی به کوهی از اطلاعات جزئی روزانه ندارند، آنها نیازمند چکیده اطلاعات برای دوره‌های زمانی متفاوت می‌باشند و به همین علت است که داده‌های تاریخی دارای مفهوم و ارزش بیشتری می‌باشد. انباره‌داده‌ها برای شرکتهایی که مصمم به استفاده از داده‌کاوی هستند یک ضرورت می‌باشد، یکی از ماهیت‌های وجودی انباره‌داده‌ها، یکپارچگی داده‌ها هنگام قرارگیری در انباره‌داده‌هاست. این بدین معنی است که وقت بسیاری به کار گرفته می‌شود تا یکنواختی و پیوستگی در درک اهداف عام سازمانی بوجود آید. اگر انباره‌داده‌ها وجود نداشته باشد، داده‌کاو باید زمان بسیار زیادی را صرف جمع‌آوری، پاکسازی و یکپارچه‌سازی داده‌ها کند. بدین ترتیب وقت بسیاری باید صرف شود تا تحلیل داده‌ها آغاز شود.

در انباره‌داده‌ها، داده‌های تاریخی جمع‌آوری و سازماندهی می‌شوند. وجود داده‌های تاریخی برای یافتن الگوها و روابطی که سازمان به دنبال آنهاست، برای داده‌کاوی یک ضرورت است. اگر چنانچه این داده‌های تاریخی وجود نداشته باشند، داده‌کاو باید به دنبال جمع‌آوری آنها باشد. علت دیگر اهمیت انباره‌داده‌ها این است که انباره‌داده‌ها شامل داده‌های جزئی و داده‌های کلی، در کنار یکدیگر می‌باشد. بدون تردید، داده‌کاو به اطلاعات جزئی برای تحلیل نیازمند است، اما داده‌های خلاصه شده نیز به کار می‌آیند. از آنجا که در انباره‌داده انواع داده‌های خلاصه شده وجود دارد، داده‌کاو می‌تواند به سرعت داده‌های انباره‌داده را بررسی کند و این باعث کاهش تکرار تحلیلها توسط داده‌کاو می‌شود.

## ۶-۲- مفاهیم انباره داده‌ها

ویژگیهای مهم یک انباره داده عبارتند از:

- موضوع محوری<sup>۱</sup>
- جامعیت<sup>۲</sup>
- پویاپذیری<sup>۳</sup> (مهم بودن عامل زمان)
- پایایی<sup>۴</sup> (غیر فرار و دائمی بودن)

**موضوع محوری:** داده‌ها طبق یک موضوع خاص سازماندهی می‌شوند، به عنوان مثال داده‌های مربوط به مشتریان، محصولات و یا داده‌های مرتبط با فروش، هر کدام جداگانه در نظر گرفته می‌شوند. اما در پایگاه داده‌های معمولی، داده‌ها بر اساس عملیات و پردازش‌های روزانه ایجاد می‌شوند و موضوع آنها مرتبط با کل پردازش می‌باشد.

**جامعیت:** داده‌های انباره، از تجمع دیگر داده‌ها ساخته می‌شوند. این داده‌ها ممکن است مربوط به پایگاه داده‌های رابطه‌ای، فایلهای بدون ساختار و یا رکوردهای مرتبط با پردازش‌های برخط باشند. جهت یکسان کردن داده‌ها، روش‌های پاکسازی به کار برده می‌شود. نوعی هماهنگی کلی در مورد داده‌های مختلفی که از سیستمهای متفاوت آمده‌اند، لازم است. به عنوان نمونه لازم است مقادیر عددی مربوط به «قیمت هتل»، «هزینه صبحانه»، «مالیات» و دیگر موارد مشابه که ممکن است از مکانهای متفاوت آمده باشند، یکسان شده و یک انباره داده با نام «مسافر» ایجاد شود.

**پویاپذیری:** افق زمانی برای انباره داده‌ها بسیار مهم‌تر از داده‌های مرتبط با سیستمهای عملیاتی می‌باشد. در ساختار انباره داده‌ها عاملی به نام زمان در نظر گرفته می‌شود. این عامل می‌تواند به‌طور ضمنی و یا به وضوح بیان شود. اما در سیستمهای عملیاتی عامل زمان، عاملی کلیدی نیست.

<sup>۱</sup>- Subject Oriented

<sup>۲</sup>- Integrated

<sup>۳</sup>- Time Variant

<sup>۴</sup>- NON Volatile

پایگاه داده‌ها شامل داده‌هایی است که روزانه با آنها کار می‌شود و بخش‌هایی به آن اضافه و یا از آن حذف می‌شود، در مقابل انباره‌داده این ویژگی را ندارد. با توجه به همین امر واضح است که بهروز شدن داده‌ها در انباره‌داده‌ها مقدور نمی‌باشد انباره‌داده‌ها نیازی به پردازش‌هایی از قبیل: تراکنشهای داده‌ای، بازیافت و مکانیزم‌های کنترل هم‌زمان ندارد. تنها اعمالی که در انباره‌داده‌ها صورت می‌پذیرد عبارتند از: بارگذاری (مقدار دهنده) اولیه داده‌ها و دسترسی به داده‌ها.

پردازشی که بر روی انباره‌داده‌ها انجام می‌گیرد *OLAP* نامیده می‌شود (*OLAP* در مقابل *OLTP* آمده است که پردازش‌هایی است که بر روی داده‌های پایگاه داده‌ها انجام می‌گیرد). جدول (۱-۶)<sup>۱</sup> و <sup>۲</sup>*OLTP* (۱-۶) مقایسه پردازشها در انباره‌داده‌ها و پایگاه داده‌ها

جدول (۱-۶) مقایسه پردازشها در انباره‌داده‌ها و پایگاه داده‌ها

<i>OLTP</i>	<i>OLAP</i>	ویژگیها
اپراتورها	کارشناسان خبره	کاربران
کارهای روزمره	تصمیم گیری	کارکرد
بر مبنای کاربرد	بر مبنای موضوع	طراحی پایگاه داده
جاری، روزبه روز و با جزئیات	تاریخی، چند بعدی مجتمع	داده
نکراری	در موارد خاص	کاربرد
خواندن و نوشتمن	کاوش و کشف	نحوه دسترسی
پردازش ساده	جستجوهای پیچیده	واحد کاری
دهها	میلیونها	تعداد رکوردها
هزاران	صدها	تعداد کاربران
۱۰۰ MB _ GB	۱۰۰ TB _ GB	اندازه پایگاه داده
جستجو و پاسخ	پردازش	شاخص

<sup>۱</sup>- Online Analytical Process

<sup>۲</sup>- Online Transaction Process

این پارامترها عبارتند از: «کاربران»، «کارکرد»، «طراحی پایگاه داده»، «داده»، «کاربرد»، «نحوه دسترسی»، «واحد کاری»، «تعداد رکوردهای در دسترس»، «تعداد کاربران»، «اندازه پایگاه داده» و «شاخص».

### ساختار انباره داده

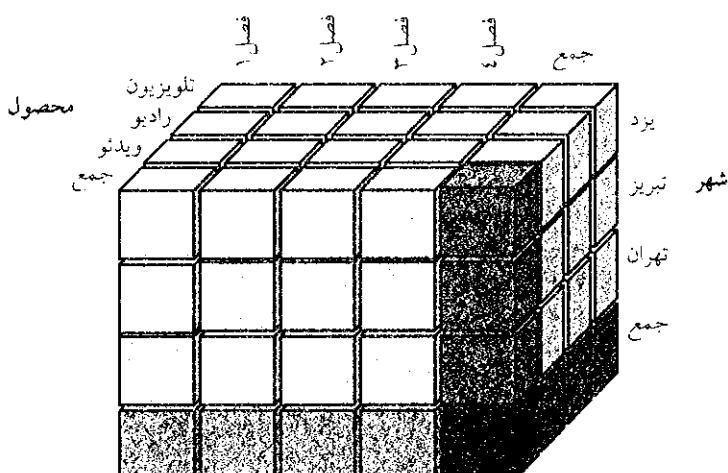
انباره دادهها بر مبنای ساختاری چند بعدی، که مدل دادهها را به شکل مکعب اطلاعاتی نشان می‌دهد. ساخته شده است. یک مکعب اطلاعاتی دادهها را در چندین بعد مختلف مدلسازی می‌کند. در زیر مثالی از یک مکعب اطلاعاتی به نام «فروش» بررسی شده است که شامل این ابعاد می‌باشد:

*Item* (نوع، نام تجاری، نام کالا)

*Time* (سال، فصل، ماه، هفته، روز)

یک مکعب اطلاعاتی اجازه می‌دهد که داده‌ها در ابعاد مختلفی مدلسازی شده و استفاده شوند. عموماً سازمانها با توجه به نیازهای آینده‌شان ابعاد مختلفی از داده‌ها را نگهداری می‌کنند. به عنوان مثال داده‌های فروش در سازمانها نیاز است با دیدگاههای زیر ذخیره شوند: در فرایند فروش مهم است که دقیقاً چه اقلامی فروخته شده‌اند. نام آنها چه بوده است نام تجاری آنها چیست و دیگر اطلاعات مرتبط. زمانهای فروش نیز مهم است.

تاریخ



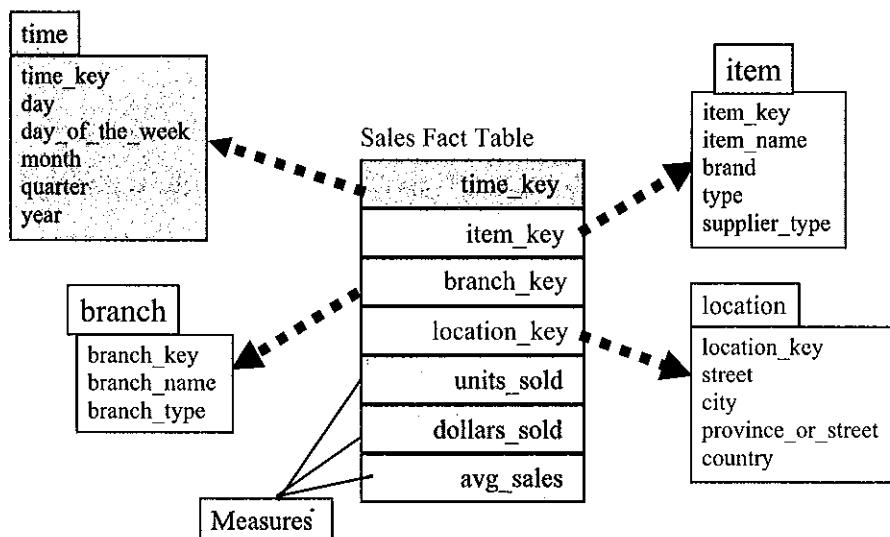
شکل ۶-۱) ابعاد داده فروش

اینکه فروش روزانه، ماهانه، فصلی و سالانه هر کدام چه میزان بوده است. اگر فروش در مکانهای مختلف جغرافیایی صورت گرفته است، هر کدام چگونه بوده‌اند. شکل (۱-۶) این ابعاد را نشان می‌دهد.

### ۶-۳-۱- مدل مفهومی انباره‌داده‌ها

در ادامه به چند مفهوم اساسی در انباره‌داده‌ها اشاره می‌کنیم:

مدل ستاره‌ای<sup>۱</sup>: یک جدول اصلی که متصل به مجموعه‌ای از جداول دیگر باشد مدل ستاره‌ای نامیده می‌شود. شکل (۶-۶) یک مدل ستاره‌ای می‌باشد که جدول اصلی فروش<sup>۲</sup> را به جداول دیگر از جمله اقلام فروش<sup>۳</sup>، مکان<sup>۴</sup> و زمان<sup>۵</sup> و شعب فروش<sup>۶</sup> متصل می‌کند.



شکل (۶-۶) مدل ستاره‌ای

<sup>۱</sup>- Star Schema

<sup>۲</sup>- Sales

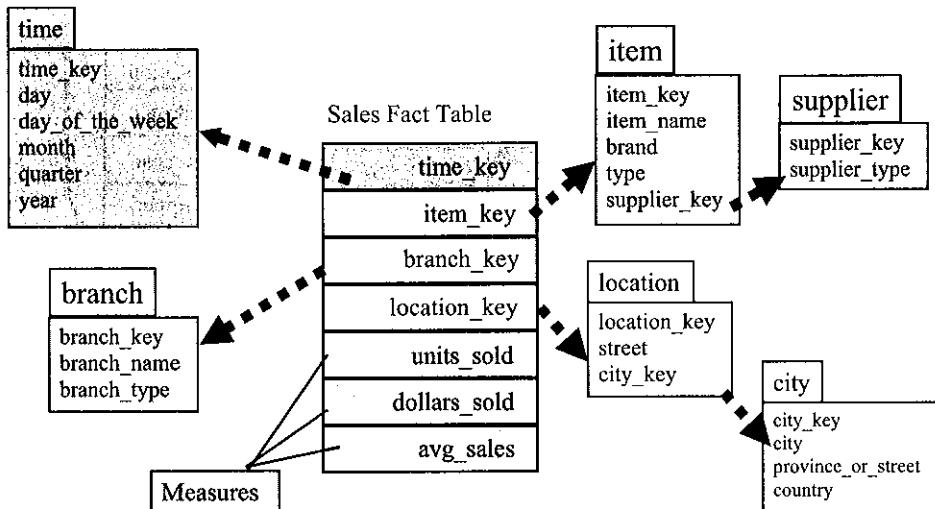
<sup>۳</sup>- Item

<sup>۴</sup>- Location

<sup>۵</sup>- Time

<sup>۶</sup>- Branch

اگر در مدل ستاره‌ای جداول جانبی با شکستن به چند جدول، نرمال شوند. این مدل تبدیل به مدل برفدانه<sup>۱</sup> می‌شود. به عنوان مثال در شکل (۳-۶) جدول مکان به دو جدول شکسته شده است و بخشی از اطلاعات آن در جدول دیگری به نام شهر<sup>۲</sup> وارد شده است.



شکل (۳-۶) مدل برفدانه‌ای

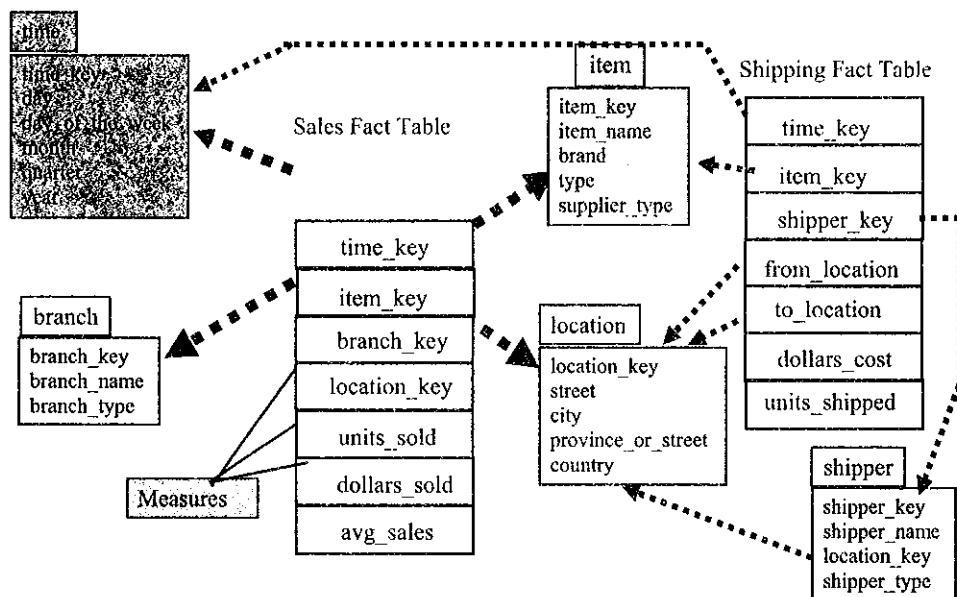
حال اگر این اطلاعات اصلی به همراه جداولش به اطلاعات دیگری نیز مرتبط باشند و بخواهیم این ارتباطات را نیز نمایش دهیم از مدل صورت فلکی<sup>۳</sup> استفاده می‌کنیم. شکل (۶-۴) ارتباطات جدول اصلی فروش را با جدول اصلی دیگری به نام حمل و نقل<sup>۴</sup> نشان می‌دهد.

<sup>۱</sup>- Snow-Flake Schema

<sup>۲</sup>- City

<sup>۳</sup>- Fact Constellation

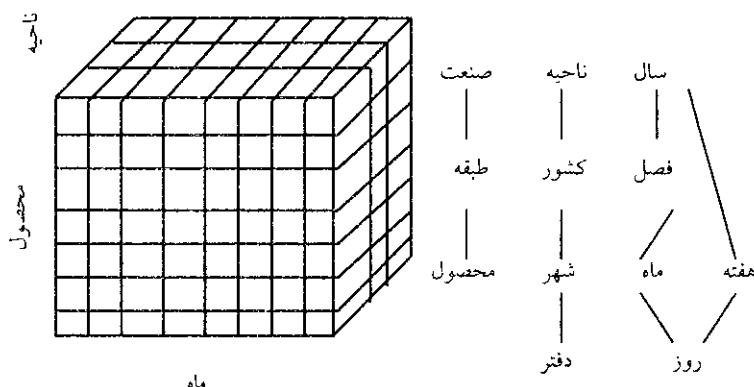
<sup>۴</sup>- Shipping



شکل ۶-۴) مدل صورت فلکی

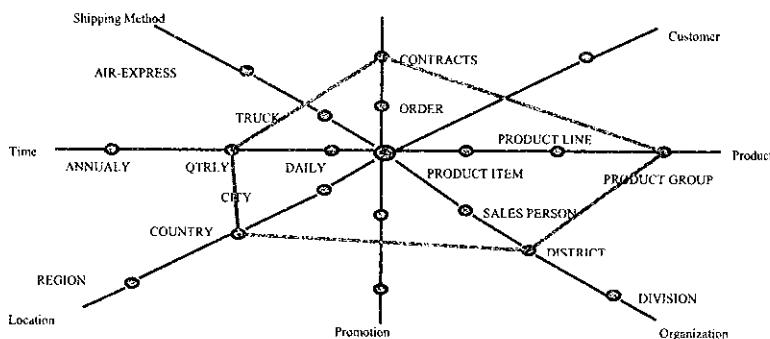
### داده‌های چند بعدی

«فروش» یک مفهوم چند بعدی است که از بخش‌های محصول، زمان فروش و محل فروش تشکیل شده است و می‌توان داده چند بعدی فروش را به شکل زیر نمایش داده در شکل (۶-۵)، سلسله مراتب ابعاد داده‌ها نمایش داده شده‌اند.



شکل ۶-۵) ابعاد مختصات داده‌های فروش

اگر ابعاد داده‌ها بیشتر از سه بعد باشد، می‌توان با مدل شبکه ستاره‌ای<sup>۱</sup> آن را نمایش داد. شکل (۶-۶) داده‌های مرتبط با ابعاد فروش را نشان می‌دهد. روی محورنماشان داده یک بعد نمایش داده می‌شود. به عنوان مثال بعد زمان روی یک محورنماشان داده شده و علاوه بر آن سلسله مراتب زمان نیز که عبارتند از روزانه، فصلی و سالانه روی این محور نشان داده می‌شوند به عنوان مثال اگر یک محصول خاص، فروش فصلی داشته و در یک مکان جغرافیایی خاص فروخته شود، در این صورت نقاط روی این نمودار که حاوی این اطلاعات می‌باشند با یک خط شکسته به هم متصل می‌شوند و بدین ترتیب داده‌های این محصول خاص با ابعاد مورد نظر نمایش داده می‌شود.



شکل ۶-۶) مدل شبکه ستاره‌ای

### زبان MQL جهت پیاده‌سازی انباره داده‌ها

زبان *MQL*<sup>۲</sup> شبیه زبان *SQL*<sup>۳</sup> می‌باشد و برای تعریف هر کدام از مفاهیم ارائه شده در بخش قبلی روش‌های خاصی وجود دارد. جهت تعریف یک جدول اصلی از دستور *Define Cube*

<sup>۱</sup>- Star Net

<sup>۲</sup>- Mining Query Language

<sup>۳</sup>- Structured Query Language

برای تعریف جداول جانبی از دستور *Define Dimension* استفاده می‌شود. مثال زیر دستورات مربوط به شکل مدل ستاره‌ای را به زبان *MQL* نشان می‌دهد.

```
define cube sales_star [time, item, branch, location]:
  dollars_sold = sum(sales_in_dollars), avg_sales = avg(sales_in_dollars) ,
  units_sold = count(*)
define dimension time as (time_key, day, day_of_week, month, quarter, year)
define dimension item as (item_key, item_name, brand, type, supplier_type)
define dimension branch as (branch_key, branch_name, branch_type)
define dimension location as (location_key, street, city, province_or_state, country)
```

دستورات زیر تعریفی از مدل برف دانه‌ای با استفاده از زبان *MQL* است.

```
define cube sales_snowflake [time, item, branch, location]:
  dollars_sold = sum(sales_in_dollars), avg_sales = avg(sales_in_dollars),
  units_sold = count(*)
define dimension time as (time_key, day, day_of_week, month, quarter, year)
define dimension item as (item_key, item_name, brand, type, supplier(supplier_key,
  supplier_type))
define dimension branch as (branch_key, branch_name, branch_type)
define dimension location as (location_key, street, city(city_key, province_or_state,
  country))
```

## ۶-۲-۳- فرایند طراحی انباره‌داده

برای طراحی یک انباره‌داده مؤثر، اولین مرحله، درک و تحلیل نیازهای کسب و کار و ساخت چارچوب تحلیل کسب و کار می‌باشد. ساخت یک سیستم اطلاعاتی پیچیده و بزرگ می‌تواند مانند ایجاد یک ساختمان بزرگ و پیچیده در نظر گرفته شود که مالک، معمار و سازنده آن دیدهای مختلفی داشته و این نظرات برای تشکیل یک چارچوب پیچیده ترکیب می‌شوند. دیدگاه‌های مختلفی در طراحی انباره‌داده وجود دارد که عبارتند از:

- **دید بالا به پایین<sup>۱</sup>**: روش بالا به پایین عبارت است از روشی که در آن ابتدا یک طرح کلی ایجاد شده و سپس به جزئیات پرداخته می‌شود.
  - **دید پایین به بالا**: ابتدا نمونه‌های کوچک ساخته شده و سپس نمونه گسترش داده می‌شود.
  - **دید پرس و جوی کسب و کار<sup>۲</sup>**: یک شکل کلی از داده‌های انباره داده بر مبنای نگرش کاربر نهایی ایجاد می‌شود.
- از دیدگاه مهندسی نرم افزار دو روش عمده جهت طراحی انباره داده‌ها وجود دارد که عبارتند از مدل آبشاری و مدل مارپیچی.
- **مدل آبشاری<sup>۳</sup>**: این مدل یک روش ساخت‌یافته و نظاممند بوده که در ایجاد انباره داده، گام به گام جلو می‌رود.
  - **مدل مارپیچی<sup>۴</sup>**: این روش یک نوع مدل‌سازی سریع است که در ابتدا یک مدل کوچک ساخته شده و سپس با بررسی مجدد آن را بهبود می‌دهند.

### ۶-۳-۲- معماری انباره داده

رویکرد چند لایه در انباره داده‌ها نیازمند این است که داده‌ها به شکلهای مختلف درآیند. این رویکرد باعث به وجود آمدن یک سیستم جامع برای مدیریت داده به منظور تصمیم‌گیری می‌شود. مهم‌ترین اجزاء این سیستم همان‌طور که در شکل (۶-۷) نشان داده شده است عبارتند از:

- سیستمهای منبع<sup>۵</sup>، جایی که داده‌ها از آنجا می‌آیند و همان سیستمهای عملیاتی می‌باشند.
- استخراج، انتقال و بارگذاری داده میان منابع مختلف داده.
- مخزن مرکزی<sup>۶</sup>، محل اصلی ذخیره‌سازی داده در انباره داده‌ها است و یک پایگاه داده رابطه‌ای با مدل منطقی می‌باشد.
- مخزن فراداده<sup>۷</sup>، توضیح می‌دهد که چه چیزهایی وجود داشته و در کجا موجود هستند.

<sup>۱</sup>- Top-Down

<sup>۲</sup>- Business Query

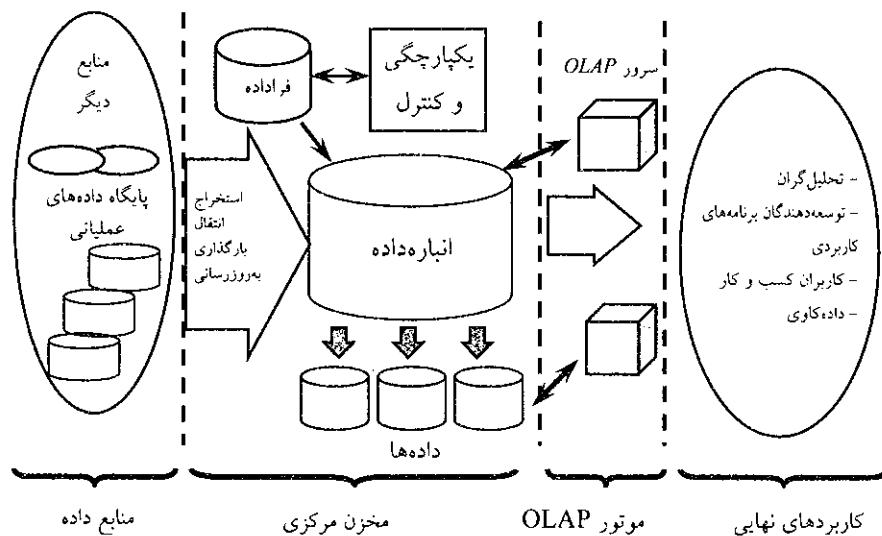
<sup>۳</sup>- Water Fall

<sup>۴</sup>- Spiral

<sup>۵</sup>- Source System

<sup>۶</sup>- Central Repository

- فرآداده، دسترسی سریع و اختصاصی را برای کاربران نهایی و برنامه‌های کاربردی فراهم می‌کند.
- بازخور عملیاتی، سیستمهای پشتیبان تصمیم را با سیستمهای عملیاتی یکپارچه می‌کند.
- کاربران نهایی، مهم‌ترین دلیل اصلی توسعه انبارهای داده در مرحله اول می‌باشند. آنها از داده‌ها و دانش استخراج شده از آنها استفاده می‌کنند.



شکل ۶-۷) معماری انباره داده ها

تقریباً همه مؤلفه هایی که ذکر شد در تمامی انباره های داده وجود دارند. داده همانند آب می باشد که از منابع سیستم سرچشم می گرفته و در انباره داده جاری شده تا به کاربران نهایی ارائه شود. این مؤلفه ها در بستر های سخت افزار، نرم افزار و شبکه سوار شده اند، این زیر ساخت ها، باید به اندازه کافی قوی باشند تا نیازمندی های کاربران نهایی و همچنین نیازمندی های پردازش و رشد داده را پوشش دهند. در ابتدا با چهار عمل اصلی استخراج<sup>۱</sup>، بهروزرسانی<sup>۲</sup>، بارگذاری<sup>۳</sup> و انتقال<sup>۴</sup>،

<sup>۱</sup>- Metadata Repository

<sup>۲</sup>- Extract

<sup>۳</sup>- Refresh

داده‌ها از پایگاه داده‌های معمولی جمع‌آوری شده و به انباره داده فرستاده می‌شوند. داده‌های موجود در انباره داده‌ها مستقیماً با دیگر سیستم‌ها در ارتباط نیستند و اگر یک سیستم عملیاتی بخواهد از داده‌های انباره داده استفاده کند از بازارچه داده‌ها<sup>۳</sup> استفاده می‌کند. بازارچه داده‌ها بخشی است که داده‌های مربوط به یک برنامه کاربردی خاص را به‌طور موقت از انباره داده‌ها دریافت کرده و در اختیار کاربر می‌گذارد. در واقع به‌جز ذخیره و بازیابی اطلاعات هیچ عملیات دیگری بر روی انباره داده‌ها امکان‌پذیر نیست.

همان‌طور که اشاره شد بازارچه داده‌ها یک سیستم تخصصی است که کلیه داده‌های مورد نیاز یک بخش یا یک برنامه کاربردی را فراهم می‌کند. بازارچه داده‌ها معمولاً در سیستم‌های گزارش‌دهی مورد استفاده قرار می‌گیرند. این قبیل بازارچه داده‌ها معمولاً از فناوری OLAP استفاده می‌کنند. نیازی نیست که تمامی اطلاعات بازارچه داده مستقیماً از مخزن مرکزی آمده باشند، در واقع مخزن مرکزی یکی از منابع داده‌ای بازارچه داده‌ها می‌باشد.

فراداده‌ها، داده‌های مربوط به داده‌ها هستند. فراداده‌ها وضعيت‌های مختلف داده‌ها را توصیف می‌کنند. مثالهایی ساده از توصیف فراداده‌ها عبارتند از:

- اطلاعات ساختاری داده‌ها چگونه ذخیره و سازماندهی شده‌اند.
- اطلاعات متريک: مقدار داده‌ها و نحوه توزيع آنها چگونه است.
- اطلاعات تجاری: داده‌ها چگونه استفاده می‌شوند.

این مخزن می‌تواند به عنوان یکی از مؤلفه‌های بانک اطلاعاتی تلقی شود. در واقع فراداده، ابزاری را در اختیار کاربران نهایی قرار می‌دهد تا به راحتی در انباره داده به جستجو بپردازند.

## ۶-۳- انواع انباره داده

از نقطه نظر معماری سه مدل انباره داده وجود دارد: انباره بنگاه<sup>۴</sup>، بازارچه داده و انباره مجازی.<sup>۵</sup>

<sup>۱</sup>- Load

<sup>۲</sup>- Transform

<sup>۳</sup>- Data Mart

<sup>۴</sup>- Enterprise Warehouse

- **انباره بنگاه اقتصادی:** این مدل کلیه اطلاعات درباره موضوعات معین داخل سازمان را گردآوری می‌کند. معمولاً از یک یا چند سیستم عملیاتی و یا ارائه کنندگان اطلاعات، داده‌ها فراهم می‌شوند. این مدل شامل داده‌های کلی و داده‌های جزء می‌باشد و می‌تواند در اندازه‌های کوچکی از گیگابایت تا هزاران گیگابایت، ترابایت و یا بیشتر مرتب شود. یک انباره داده بنگاه می‌تواند تحت سیستمهای مین‌فریم<sup>۲</sup> است، ابرسرورهای unix یا بسترهاي معماری موازی پیاده‌سازی شود. این انباره نیازمند مدلسازی کسب و کار گسترشده می‌باشد که طراحی و ساخت آن ممکن است سالهای زیادی به طول انجامد.
- **بازارچه داده:** شامل زیرمجموعه‌ای از داده‌های سازمانهای گسترشده است که شامل داده‌های مرتبط با گروه ویژه‌ای از کاربران می‌باشد. به عنوان مثال یک بازارچه داده بازاریابی می‌تواند به موضوعاتی مانند مشتری، اقلام جنس و فروش محدود شود. داده‌های موجود در بازارچه داده تمایل به خلاصه شدن دارند. بازارچه‌های داده اغلب تحت سرورهای اداری ارزان قیمت که بر مبنای windows/NT یا OS/۲ هستند، پیاده‌سازی می‌شوند. چرخه پیاده‌سازی بازارچه‌های داده بیشتر در مقیاس هفتة برآورد می‌شود. بازارچه‌های داده بر اساس منبع داده به دو دسته زیر تقسیم می‌شوند:
  - **بازارچه‌های داده مستقل:** این نوع بازارچه‌های داده مرتبط با بیش از یک سیستم عملیاتی بوده و یا مرتبط با داده‌هایی هستند که به‌طور محلی درون یک بخش خاص یا محدوده جغرافیایی خاص تولید شده‌اند.
  - **بازارچه‌های داده وابسته (غیر مستقل):** به‌طور مستقیم از انباره بنگاه استخراج می‌شوند.
- **انباره مجازی:** این انباره مجموعه‌ای از برشها بر اساس دیدگاهها<sup>۳</sup> مختلف بر روی پایگاههای داده عملیاتی می‌باشد. برای پردازش یک پرس‌وحجی مؤثر فقط برخی از

'- Virtual Warehouse

'- Main Frame

'- View

دیدگاه‌های خلاصه به کار می‌رود. ساخت یک انباره مجازی آسان است اما نیازمند ظرفیت اضافه در سرورهای مربوطه می‌باشد.

## ۶-۴- انباره داده و سیستم‌های عملیاتی

انباره داده مخزنی از داده‌های یک بنگاه است که اغلب برای تحقیق و پشتیبانی تصمیمات از آن استفاده می‌شود. این انباره در مقابل سیستم عملیاتی سازمان<sup>۱</sup> که با تراکنشهای روزمره سازمان سروکار دارد (مثلًا OLTP) قرار می‌گیرد. از آنجایی که سیستمهای عملیاتی اغلب به صورت مستقل طراحی می‌شوند، برای پردازش یکپارچه داده‌ها به مخزنی خاص نیاز دارند که کلیه داده‌های مرتبط را یکجا دربرگیرد.

مزیت دیگر انباره داده این است که فعالیتهای تحلیلی مانند OLAP را از سیستم عملیاتی جدا می‌کند. از انباره داده‌ها می‌توان برای داده‌کاوی، مصورسازی داده‌ها<sup>۲</sup>، ارائه گزارشات پیشرفته و به کارگیری انواع ابزارهای OLAP استفاده کرد. همچنین اطلاعات از منبع اصلی مستقل می‌شوند که این مسئله در زمانی که اطلاعات اجرایی در حال تغییر هستند بسیار اهمیت دارد. علاوه بر این اگر ساختار ذخیره اطلاعات جهت یک نوع عملیات خاص، طراحی شده باشد (مثلًا ساختار ستاره‌ای برای عملیات فروش) جستجو در ابعاد مختلف نیز بسیار ساده‌تر خواهد شد. با توجه به آنچه گفته شد می‌توان در موارد زیر انباره داده‌ها را از سیستمهای عملیاتی متمایز نمود:

- اهداف
- ساختار
- اندازه
- بهینه بودن عملکرد
- فناوری‌های استفاده شده

<sup>۱</sup>- Operational System

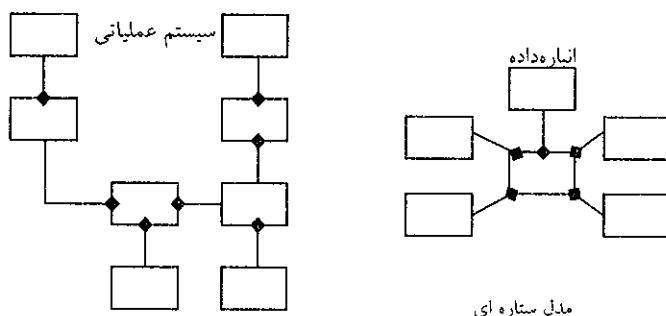
<sup>۲</sup>- Data Visualization

جدول (۶-۲) این تفاوت‌ها را به تفصیل نشان می‌دهد:

جدول (۶-۲) مقایسه سیستم‌های عملیاتی و انباره‌داده‌ها

انباره‌داده‌ها	سیستم‌های عملیاتی
بر مبنای موضوع ( <i>hundreds of GB up to TB</i> )	بر مبنای کاربرد ( <i>several MB up to GB</i> )
داده‌های تاریخی	داده‌های جاری
جداول غیرنرمال	جداول نرمال
به روز شدن دسته‌ای	به روز شدن همیشگی
جستجوهای پیچیده	جستجوهای ساده

ساختار اغلب انباره‌ها به صورت ستاره‌ای طراحی می‌شود زیرا اولاً موضوع‌گرا هستند و ثانیاً از آنجا که هر نوع ارتباطی بین موجودیتها ممکن است مورد مطالعه قرار گیرد، تا حد امکان لازم است موجودیتها بدون واسطه مرتبط شده باشند.



نمودار

شکل (۸-۶) تفاوت ساختارها

شکل (۸-۶) تفاوت در ساختار دو نوع سیستم را نشان می‌دهد. عامل مهم دیگری که در انباره‌داده‌ها باید به آن توجه شود نوع و ماهیت داده‌ها است چرا که داده‌ها باید دقیق<sup>۱</sup>، سازگار،

<sup>۱</sup>- Accurate

به موقع<sup>۱</sup>، یکپارچه، کامل<sup>۲</sup>، معتبر<sup>۳</sup> بوده و در راستای قواعد کسب و کار باشند و به علاوه باید به خوبی قابل درک<sup>۴</sup> باشند.

#### ۶-۴-۱- کاربران نهایی انباره داده‌ها

کاربران نهایی در واقع آخرين و مهم‌ترین مؤلفه در هر انباره داده می‌باشند. اين کاربران نهایي تحليل گران، توسعه دهنده‌گان برنامه‌های کاربردي و کاربران کسب و کار می‌باشند.

#### تحلیل گران

تحلیل گران نیاز دارند که به غالب داده‌ها به منظور استخراج مدل‌های مختلف و تهیه گزارشها دسترسی داشته باشند. آنها از یکسری ابزارهای خاص از جمله بسته‌های آماری، ابزارهای داده‌کاوی و صفحات گسترده استفاده می‌کنند. عموماً تحلیل گران به عنوان نخستین ذینفعان انباره‌های داده محسوب می‌شوند. تعداد افراد خبره‌ای که در اين دسته قرار می‌گيرند بسیار کم است. کاري که آنها انجام می‌دهند از درجه اهمیت بسیار بالايی برخوردار بوده و بسیار پیچیده است. يك انباره داده، داده‌های پاکسازی شده را به طور يكجا جمع‌آوری می‌کند. اين داده‌ها باید ویژگیهای زیر را دارا باشند تا بتوانند به راحتی مشکلات تحلیل گران را حل کنند:

- داده‌های سراسر بانک اطلاعاتی باید سازگار باشند.
- داده‌ها باید با زمان سازگار باشند.
- يك سистем باید بتواند به پایین‌ترین سطح اطلاعات و تراکنشها دسترسی داشته باشد.

#### توسعه دهنده‌گان برنامه‌های کاربردی

انباره‌های داده عموماً طیف گسترده‌ای از برنامه‌های کاربردی را پشتیبانی می‌کنند. به منظور توسعه يك برنامه کاربردی پايدار انباره‌های داده نقش بهسزايی دارند. اول اينکه برنامه‌ای را که آنها توسعه می‌دهند باید در برابر تغييرات در ساختار انباره داده ايمان باشد. ايجاد جداول جديد،

<sup>۱</sup>- Consistent

<sup>۲</sup>- Timely

<sup>۳</sup>- Complete

<sup>۴</sup>- Valid

<sup>۵</sup>- Well Understood

فیلدهای جدید و تغییرات ساختاری جداول باید حداقل تأثیر را بر روی برنامه‌های کاربردی موجود داشته باشد. وجود یکسری مشخصه ویژه بر روی داده‌ها به تحقق این امر کمک می‌کند. علاوه بر آن داشتن دانشی درباره اینکه هر برنامه کاربردی از چه فیلدهایی استفاده می‌کند، می‌تواند مانع بن بست<sup>۱</sup> شود. توسعه دهنده‌گان سیستم نیاز دارند بدانند ارزش معتبر فیلدها چیست و علاوه بر آن ارزش هر فیلد به چه معناست. پاسخ این سوالات هدف فراداده می‌باشد. فراداده مستنداتی را در ارتباط با ساختار داده ارائه می‌کند. از آنجا که کسب و کار واقعی نیازمند توسعه برنامه‌های کاربردی می‌باشد، درک نیاز توسعه دهنده‌گان و دسترسی آنها به انباره‌داده‌ها از درجه اهمیت بالایی برخوردار است. انباره‌های داده به مرور چهار تغییر می‌شوند و برنامه‌های کاربردی نیز همچنان از آنها استفاده می‌کنند. این امر مهم‌ترین عامل موفقیت کنترل و مدیریت تغییرات می‌باشد.

### کاربران کسب و کار

کاربران کسب و کار نیز از جمله استفاده‌کنندگان انباره‌داده می‌باشند. نیازمندیهای آنها موجب توسعه برنامه‌های کاربردی و معماری انباره‌داده می‌شود. در برخی از کسب و کارها کاربران نهایی تنها با گزارش‌های تهیه شده از انباره‌های داده، یا صفحات گسترده سروکار دارند. اینکه افراد بر روی میز خود کامپیوتر داشته و قادر باشند به طور مستقیم به انباره‌داده دسترسی داشته باشند، از درجه اهمیت بالایی برخوردار است. کاربران با استفاده از ابزارهای موجود می‌توانند گزارش‌های ترسیمی و تحلیلی بسیار جالب و با ارزشی مناسب با نیاز خود از انباره‌داده استخراج کنند. از طرف دیگر آنها همچنین قادر خواهند بود تا به درون مخزن انباره‌داده وارد شده و تا پایین ترین سطح، داده‌های موجود را بررسی کنند.

### کاربردهای انباره‌داده

سه کاربرد مهم برای انباره‌داده شناخته شده است: یکی از این کاربردها، داده‌کاوی می‌باشد که این ارتباط قبلاً به تفصیل توضیح داده شد. کاربردهای دیگر انباره‌داده در پردازش اطلاعات و

پردازش‌های تحلیلی می‌باشد. به عنوان مثال پرس‌وجوها و تحلیلهای آماری و ایجاد جداول، نمودارها و گزارش‌های گرافیکی با استفاده از پردازش اطلاعات در انباره‌داده‌ها قابل ارائه می‌باشد. انباره‌داده‌ها به تحلیل گران کسب و کار کمکهای شایانی می‌کند:

- داشتن انباره‌داده، یک مزیت رقابتی است چرا که با دسترسی سریع و به موقع به داده‌ها به غلبه بر رقبا کمک می‌کند.
- یک انباره‌داده می‌تواند بهره‌وری کسب و کار را توسعه دهد. چرا که قادر است اطلاعات سازمان را خیلی دقیق، با سرعت و به‌طور کارا تشریح کند.
- یک انباره‌داده بازاریابی، ارتباط با مشتریان را تسهیل می‌کند، چرا که داده‌های مرتبط با مشتریان و اقلام مختلف در کلیه بخش‌ها و کلیه فروشگاهها را در اختیار قرار می‌دهد. همچنین انباره‌داده‌ها می‌توانند توسط ردیابی روندها، الگوهای استثنایات در دوره‌های زمانی طولانی با یک روش سازگار و قابل اطمینان به کاهش هزینه‌ها کمک کنند.

## منابع

- 1) Han J. and Kamber M. (2006) *Data Mining: Concepts and Techniques* San Francisco , CA: Morgan Kaufmann.

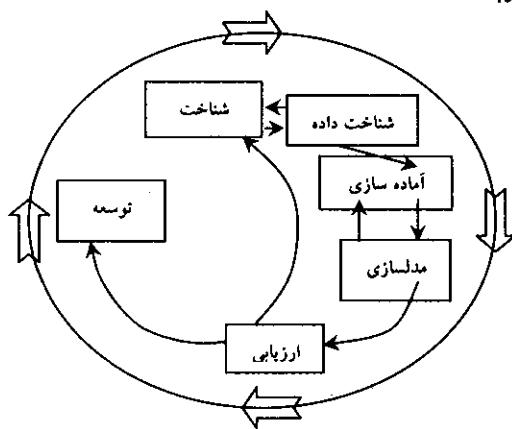
---

## فصل هفتم

---

# متداولوژی اجرا و پیاده‌سازی پروژه‌های داده‌کاوی

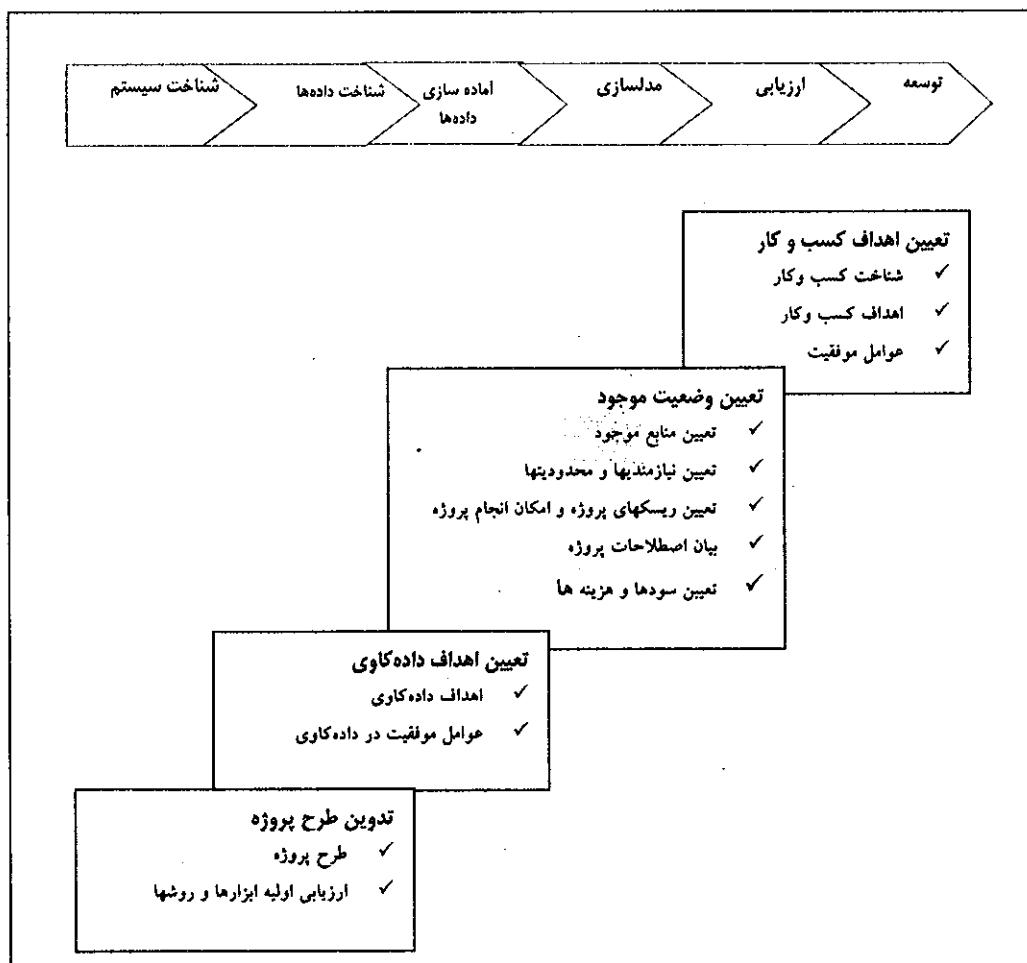
روشهای مختلفی برای پیاده‌سازی و اجرای پروژه‌های داده‌کاوی وجود دارد. یکی از روش‌های بسیار قوی، متداولوژی CRISP<sup>۱</sup> می‌باشد. این متداولوژی از گامهای شناخت سیستم، شناخت داده‌ها، آماده‌سازی داده‌ها، مدل‌سازی، ارزیابی و توسعه سیستم تشکیل شده است. هر کدام از این گامها به زیر بخش‌هایی تقسیم می‌شوند.



شکل ۷-۷) گامهای متداولوژی CRISP

## ۱-۷- گام شناخت سیستم

در گام شناخت سیستم ابتدا به شناخت کسب و کار مورد نظر پرداخته می‌شود. سپس اهداف مورد نظر و عوامل موفقیت کلیدی آن تعیین شده و دوباره اهداف کسب و کار بازنگری می‌شود. شناسایی فرصتها و عوامل موفقیت کلیدی یک کسب و کار قدم بسیار مهمی می‌باشد، چرا که باعث افزایش اطلاعات و انجام بهتر کارها می‌شود. در واقع در این گام به تعیین زمینه‌ها و عواملی می‌پردازیم که داده‌ها باعث افزایش ارزش در آنها می‌شوند.



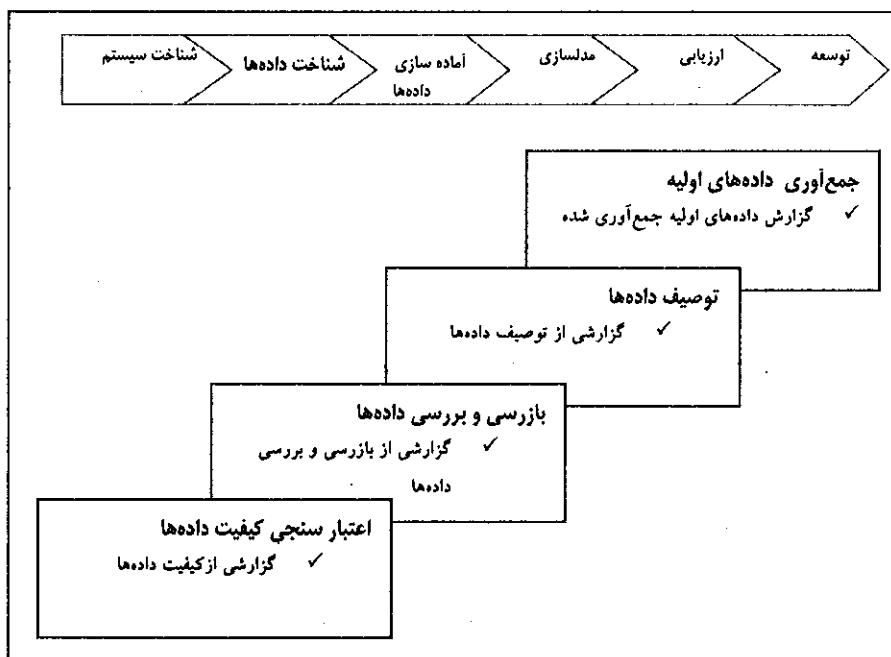
شکل ۲-۷) جزئیات مربوط به گام شناخت سیستم

پس از تدوین اهداف مورد نظر کسب و کار، می‌بایست به شناخت وضعیت موجود پرداخت. به منظور تعیین دقیق وضعیت موجود، به تعیین منابع موجود پرداخته و نیازمندیها و محدودیتهای موجود تعیین می‌شوند. شناخت ریسکهای بر سر راه پروژه و نیازمندیهای پروژه کمک می‌کنند تا طرح امکان‌سنگی پروژه تدوین شود. در این قدم به منظور تدوین طرح امکان‌سنگی، تعیین منابع پروژه ضروری می‌باشد. منابع پروژه عبارتند از: منابع انسانی، منابع مالی، تجهیزات و دیگر منابع. به همین دلیل در این گام تعیین سودها و هزینه‌های پروژه امری اجتناب‌ناپذیر می‌باشد. تعیین اهداف داده‌کاوی و تدوین طرح پروژه نیز از دیگر بخش‌هایی است که در این گام باید به آنها پرداخت. به عنوان مثال در یک کسب و کار خاص اهداف زیر برای پروژه داده‌کاوی تعیین می‌شود:

- برنامه‌ریزی بازاریابی محصولات و خدمات جدید
- قیمت‌گذاری محصولات و خدمات جدید
- شناخت مشتریان ناراضی و جلوگیری از رویگردانی آنها
- تعیین بازار هدف

## ۷-۲- گام شناخت داده‌ها

پس از شناخت کسب و کار به سراغ شناخت داده‌ها می‌رویم. شناخت داده‌ها عبارت است از جمع‌آوری داده‌های اولیه، توصیف داده‌ها، بازرگانی و بررسی داده‌ها و اعتبار‌سنگی کیفیت داده‌ها. کارآیی داده‌کاوی مستقیماً مرتبط با داده‌های مورد استفاده دارد. هر اندازه داده‌ها دقیق‌تر جامع‌تر و با کیفیت بهتری باشند خروجی داده‌کاوی کارآتر خواهد بود. بنابراین انتخاب و جمع‌آوری داده‌های درست، توصیف آنها، یکپارچه‌سازی قالب آنها به منظور استفاده در داده‌کاوی، از اهمیت بسیار بالایی برخوردار می‌باشد. علاوه بر این بازرگانی و بررسی داده‌ها به منظور تعیین میزان کیفیت آنها بسیار مهم می‌باشد. در شکل (۳-۷) به این گامها و خروجی‌های هر مرحله اشاره شده است. پروژه داده‌کاوی باید در راستای اهداف داده‌کاوی صورت گیرد. بنابراین جمع‌آوری و یکپارچه کردن داده‌های مرتبط با این اهداف در برآورده کردن اهداف و شاخصهای مرتبط با آن بسیار مهم و حیاتی است.

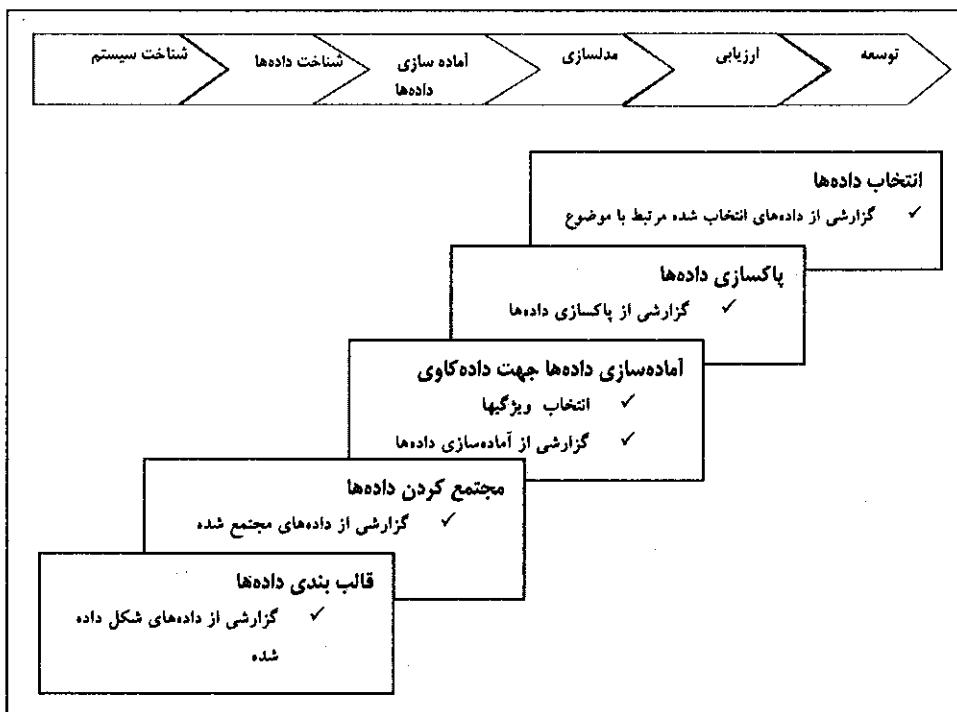


شکل ۳-۷) جزئیات مربوط به گام شناخت داده‌ها

### ۳-۷- گام آماده سازی داده‌ها

گام آماده سازی داده‌ها عبارت است از: انتخاب داده‌ها، پاکسازی داده‌ها، آماده کردن داده‌ها جهت داده کاوی، مجتماع کردن آنها و قالب بندی داده‌ها. برای اجرای هر کدام از این زیر بخشها، معالیهای دیگری نیز ضروری است که در شکل (۴-۷) آمده است.

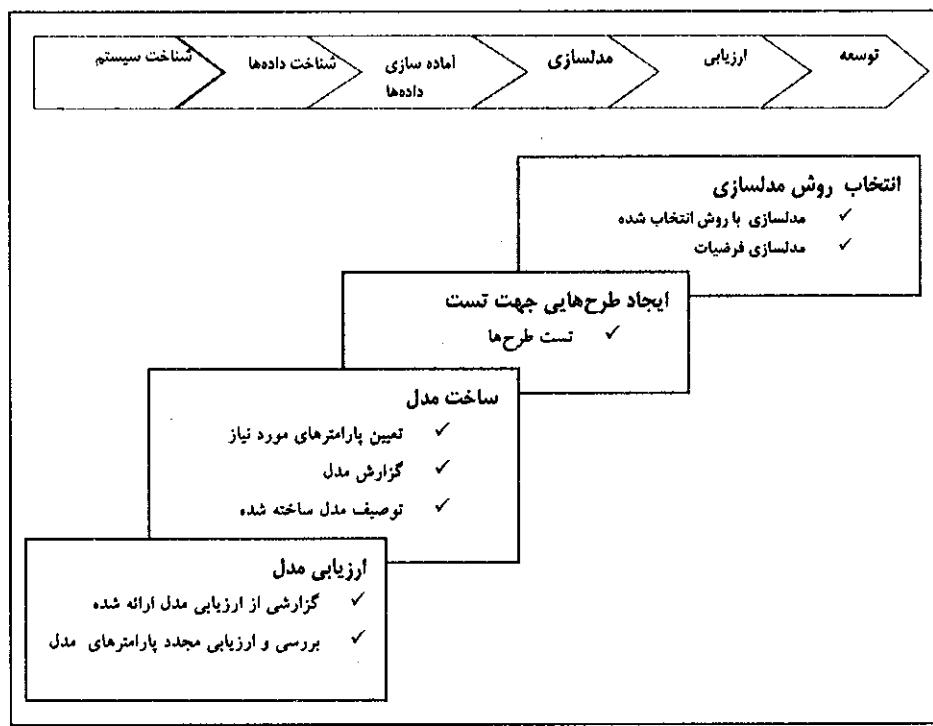
جمع اورن ر محافظت از داده‌ها گام بسیار مهمی می‌باشد. اصولاً چون قالب و نوع داده‌ها در طول زمان تغییر می‌کند، ممکن است قالب بسیاری از داده‌های موجود متفاوت باشد. همچنین به علت اینکه داده‌ها از منابع مختلف داخلی و خارجی جمع آوری شده و یکپارچه می‌شوند، باز هم ممکن است قالب داده‌ها با هم یکسان نبوده و یا حتی برخی از داده‌های قبلی از بین رفته و دور ریخته شده باشند و بخشهایی از داده‌ها موجود باشد. در داده کاوی اهمیت داده‌های قدیمی به هیچ وجه کمتر از داده‌های جدید نمی‌باشد.



شکل ۷-۴) جزئیات مربوط به گام آماده‌سازی داده‌ها

#### ۷-۴- گام مدلسازی

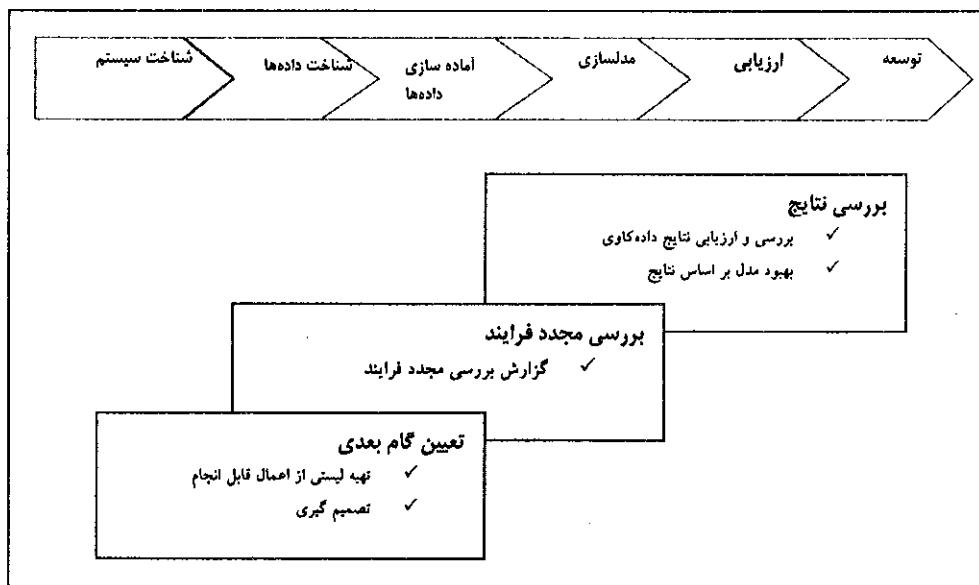
پس از شناخت داده‌ها و آماده‌سازی آنها، حال می‌توان به مدلسازی پرداخت. در اولین قدم از مدلسازی می‌باشد روش مناسب را انتخاب کرد. انتخاب روش مناسب بسیار تعیین کننده می‌باشد. پارامترهای مورد نیاز مدل نیز پس از تعیین روش مورد استفاده، مشخص می‌شوند. پس از انتخاب مدل و تعیین پارامترها، بخش‌های کوچکی از پروژه تعریف شده و پس از اجرا شدن، در هر مرحله به دقت تست می‌شوند تا کیفیت مدل ایجاد شده تضمین شود. در این مرحله اگر مدل مورد نظر دقت لازم را نداشت و یا کیفیت مطلوب را حاصل نکرد، ابتدا به تغییر پارامترهای مدل می‌پردازیم و مجدداً مدل را تست می‌کنیم. اگر هنوز کیفیت لازم را کسب نکرده بود، مدل را تغییر داده و مدل جدیدی می‌سازیم. برای انجام هر کدام از زیر بخش‌های مدلسازی، فعالیتهای دیگری نیز ضروری است که در شکل (۵-۷) آمده است.



شکل ۵-۷) جزئیات مربوط به گام مدلسازی

### ۵-۷- گام ارزیابی

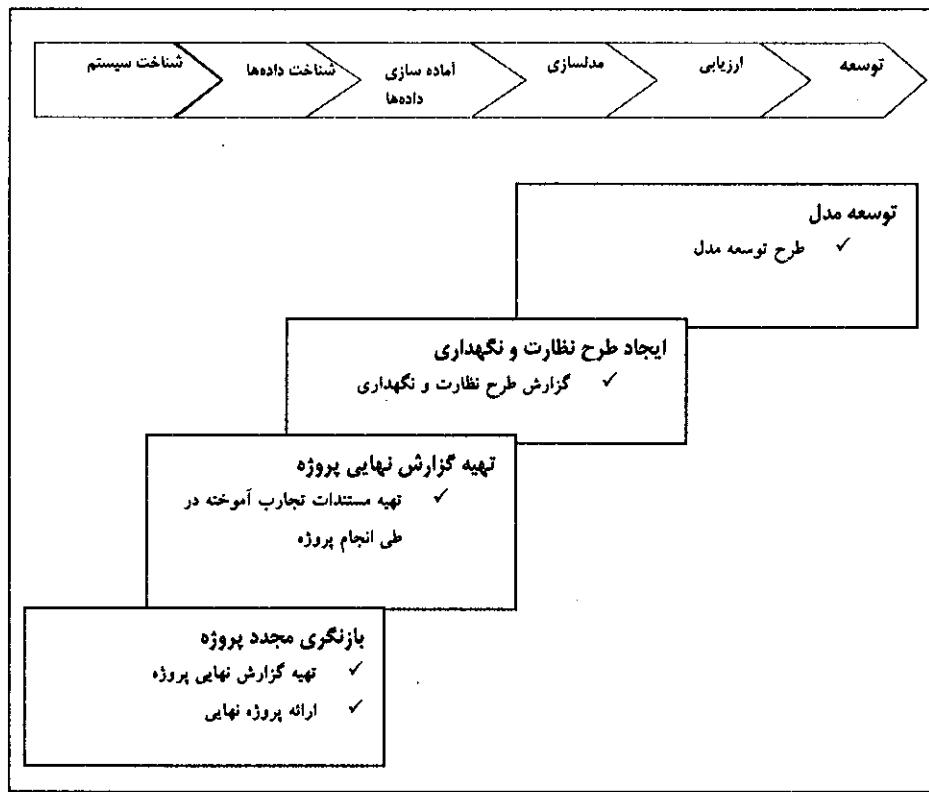
پس از مدلسازی، حال می‌بایست به ارزیابی نتایج حاصل از مدل پرداخت. نتایج ارزیابی باعث بهبود مدل شده و مدل را قابل استفاده می‌کند. در این گام اعتبار مدل بررسی شده و گزارشی از کل فرایند تهیه می‌شود. در انتهای نیز لیستی از اقدامات اصلاحی قابل انجام تهیه شده و بد عنوان راهکار ارائه شده و تصمیم‌گیریها بر این اساس انجام می‌شود.



شکل ۷-۷) جزئیات مربوط به گام ارزیابی

## ۷-۶- گام توسعه

پس از استخراج لیست اقدامات قابل انجام، دورنمایی از طرح توسعه ایجاد می‌شود. در این گام این طرح مجدداً بررسی شده و مدون می‌شود. علاوه بر آن طرح نظارت و نگهداری پس از اتمام پروژه نیز در این گام تهیه می‌شود. در پروژه‌های داده‌کاوی تأکید زیادی بر روی مستندسازی تجارب آموخته در طی انجام پروژه وجود دارد و در این گام گزارش مرتبط با آن نهایی می‌شود. در این مرحله پروژه قابل ارائه نهایی بوده و گزارش نهایی آن نیز استخراج شده است.



شکل ۷-۷) جزئیات مربوط به گام توسعه

## منابع

راهنمای نرم افزار SPSS Clementine



---

---

## بخش چهارم

---

### مباحث ویژه در داده کاوی

فصل هشتم: سریهای زمانی در داده کاوی

فصل نهم: شبکه‌های اجتماعی

فصل دهم: کاربرد داده کاوی در مدیریت ارتباط با مشتری



---

## فصل هشتم

---

# سریهای زمانی در داده‌کاوی

یک سری زمانی دنباله‌ای از مشاهدات بر روی یک متغیر مورد توجه است که در نقاط گستره‌ای از زمان که معمولاً فاصله‌های مساوی دارند (روزانه - هفتگی - ماهانه - فصلی - سالانه)، رخ می‌دهد. تجزیه و تحلیل سریهای زمانی، مخصوصاً توصیف فرآیند یا پدیده‌ای است که تولید دنباله می‌کند. جهت پیش‌بینی سریهای زمانی، لازم است که رفتار فرآیند را با یک مدل ریاضی که قابل تعمیم به آینده باشد، توصیف کرد. معمولاً لازم نیست مدل نماینده مشاهدات خیلی قدیمی یا فراتر از زمان مورد انتظار پیش‌بینی باشد<sup>[۱]</sup>. مهم‌ترین نکته در داده‌های سری زمانی، آن است که این داده‌ها دارای همبستگی هستند. از ضریب همبستگی برای تعیین همبستگی بین مقادیر  $X$  و  $Y$  استفاده می‌شود، اما وقتی خود متغیرهای مستقل، مقادیرشان به هم مرتبط باشد، به آن خودهمبستگی<sup>۱</sup> گویند. با توجه به تعریف داده‌های سری زمانی، روش‌های آماری مبتنی بر فرض مستقل بودن مشاهدات، مناسب نبوده و به جای آن می‌توان از معادلات خودهمبستگی در تحلیل سریهای زمانی استفاده کرد.<sup>[۲]</sup>

با توجه به اهمیت و نقش سریهای زمانی، در حوزه‌های مختلفی از آنها استفاده می‌شود که نمونه‌هایی از کاربردهای آن به شرح زیر است.<sup>[۲]</sup>

- بازرگانی و اقتصاد: مقادیر فروش و قیمت‌های ماهیانه، قیمت سهام در روزهای مختلف

---

<sup>۱</sup>- Auto Correlation

- علوم مهندسی: ثبت علائم الکتریکی، ولتاژ و سیگنالها.
- پژوهشکی: داده‌های الکتروکاربودیوگرام و دیگر کاربردها.
- هواشناسی: درجه حرارت روزانه و میزان بارندگی سالیانه.
- کنترل کیفیت: ثبت مشاهدات مربوط به فرایند در نمودارهای کنترل.
- علوم اجتماعی: نرخهای زاد و ولد و مرگ و میر سالیانه.

## ۱-۸- داده‌کاوی سریهای زمانی<sup>۱</sup>

یک سری زمانی ساده‌ترین شکل داده‌های زمانی است. سری زمانی دنباله‌ای از اعداد حقیقی است که به صورت منظم در طول زمان گردآوری شده است. هر عدد، نشان‌دهنده مقدار یک متغیر مشاهده شده می‌باشد. همان‌طور که اشاره شد، داده‌های سری زمانی در حوزه‌های مختلفی مثل تحلیل بازار سهام، علوم ارتباطات، پژوهشکی، داده‌های مالی و غیره مطرح می‌شوند. همچنین داده‌های وب که میزان استفاده از وب سایتها مختلف را ثبت می‌کنند (برای مثال تعداد کلیکها) را می‌توان با سریهای زمانی مدل کرد. در حقیقت، سریهای زمانی برای نمایش بخش بزرگی از داده‌های ذخیره شده در بانکهای اطلاعاتی تجاری به کار می‌رود که به تدریج به عنوان یک نوع داده متفاوت، اهمیت بیشتری یافته است. اهمیت داده‌های سریهای زمانی موجب تحقیقات زیادی در زمینه تحلیل این نوع داده‌ها شده است. ادبیات آماری در مورد سریهای زمانی بسیار وسیع است و به طور عمده به مسائلی مانند شناسایی الگوها و تحلیل روند (مانند رشد خطی فروش شرکت در طول یک سال)، تحلیلهای فصلی (مثلًا فروش زمستانی یک محصول تقریباً دو برابر فروش تابستانی است) و پیش‌بینی (مانند پیش‌بینی فروش فصل آینده) می‌پردازند. این موضوعات کلاسیک در تعدادی از مراجع آماری بررسی شده است. [۵]

هدف اصلی داده‌کاوی سریهای زمانی کشف الگوهای موجود در واقعی و داده‌های یک سری زمانی است. کشف این الگوهای ناشناخته، همگام با استفاده از دیگر روش‌های مختلف داده‌کاوی مانند سیستمهای پایگاه داده، آمار، یادگیری ماشینی، شبکه‌های عصبی، تئوری مجموعه‌ها، منطق فازی و غیره در داده‌کاوی اتفاق می‌افتد. از مهم‌ترین کاربردهای داده‌کاوی سریهای زمانی،

می‌توان به دسته‌بندی، خوشبندی و کشف قواعد از داده‌ها اشاره کرد. در داده‌کاوی سریهای زمانی دو سؤال اساسی زیر مطرح می‌شود:

- چگونه می‌توان روابط همبستگی در درون سریهای زمانی را پیدا کرد؟
  - چگونه می‌توان سریهای زمانی با حجم انبوهی از داده‌ها را تحلیل کرده و الگوهای منظم، روند، تغییرات تصادفی، داده‌های مغلوش و غیره را از آنها استخراج کرد؟<sup>[۶]</sup>
- در ادامه به بررسی جنبه‌های مختلفی از داده‌کاوی سریهای زمانی، با تمرکز بر شناسایی اجزاء سریهای زمانی و روشهای جستجوی تشابه در سریهای زمانی پرداخته می‌شود.

### ۸-۱-۱- اجزاء سریهای زمانی و تحلیل آنها

یک سری زمانی شامل یک متغیر وابسته (u) می‌باشد که تابعی از زمان است. چنین تابعی به شکل یک نمودار سری زمانی نمایش داده می‌شود. در مطالعه داده‌های سریهای زمانی، همواره دو هدف اساسی زیر دنبال می‌شود:<sup>[۶]</sup>

- مدلسازی سریهای زمانی با تأکید بر فرآیند ایجاد سریهای زمانی
  - پیش‌بینی سریهای زمانی با تأکید بر پیش‌بینی مقادیر متغیرهای سری زمانی تجزیه و تحلیل روند شامل شناسایی چهار جزء یا مشخصه اساسی هر سری زمانی می‌باشد.<sup>[۲]</sup>
  - روند خطی و غیرخطی: تغییر دراز مدت در میانگین یا به عبارت دیگر حرکت دراز مدت تدریجی افزایشی یا کاهشی داده‌ها در طول زمان است.
  - تغییر سیکلی: تغییرات دوره‌ای موجود در داده‌های یک سری زمانی است. معمولاً این نوع افزایشها و کاهشها لحظه‌ای در داده‌ها، در دوره‌های بیشتر از یکسال اتفاق می‌افتد.
  - تغییرات فصلی: تغییراتی است که به صورت فصلی در داده‌های سری زمانی اتفاق می‌افتد. این نوع الگوی تغییر در طول یکسال مشاهده می‌شود.
  - تغییرات باقیمانده‌ها یا تصادفی: اگر سه جزء قبلی از یک سری زمانی حذف شود، سری باقیمانده حاصل می‌شود. که ممکن است تصادفی باشد.
- با توجه به توضیحات ارائه شده، می‌توان معادله یک سری زمانی را به صورت یکی از دو حالت (۱-۸) و (۲-۸) نشان داد. <sup>[۶]</sup>

$$y = T_i + C_i + S_i + R_i \quad (1-8)$$

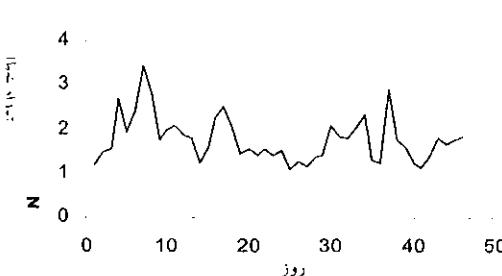
$$y = T_i * C_i * S_i * R_i \quad (2-8)$$

یکی از مهم‌ترین جنبه‌های کاربرد سریهای زمانی، پیش‌بینی می‌باشد. پیش‌بینی سریهای زمانی با استفاده از یک معادله ریاضی، الگویی تاریخی در داده‌های سری زمانی ایجاد می‌کند. این روش برای پیش‌بینی کوتاه مدت یا بلند مدت مقادیر آینده مورد استفاده قرار می‌گیرد. روش‌های مختلفی برای پیش‌بینی سریهای زمانی، مورد استفاده است که از بین آنها روش میانگین متحرک تلفیق شده با اتو رگرسیون<sup>۱</sup> که به مدل «باکس- جنکینز» نیز موسوم است از اهمیت ویژه‌ای برخوردار است[۶]. برای آشنایی بیشتر با روش‌های پیش‌بینی سریهای زمانی، می‌توان به مراجع آمار و اقتصادستنجی مراجعه کرد.

معمولًا مطلوب است که سیستم پیش‌بینی بتواند تغییرات پایدار را مشخص و با تعديل مدل پیش‌بینی، فرآیند جدید را تعقیب کند. در عین حال سیستم پیش‌بینی، تغییرات تصادفی و موقتی را تشخیص داده و در مقابل آنها واکنش نشان ندهد. در هنگام کار با مدل‌های پیش‌بینی با توجه به مراحل مختلف سیکل پیش‌بینی (مثلًا عمر محصول)، لازم است که مدل‌های پیش‌بینی مختلفی به کار گرفته شوند. مثلاً ممکن است گاهی اوقات فقط روند را حفظ کرده و بقیه علل تغییرات سری زمانی را حذف کنیم. [۴]

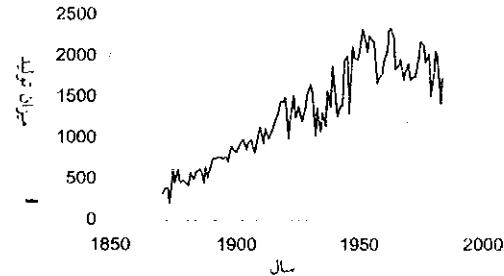
به طور کلی یک سری زمانی را ایستا<sup>۲</sup> گویند، هرگاه تغییر منظمی در میانگین و واریانس آن وجود نداشته و تغییرات دوره‌ای اکید حذف شده باشد. نظریه احتمال سریهای زمانی بیشتر با سریهای زمانی ایستا سر و کار دارد و به این دلیل است که در تجزیه و تحلیل سریهای زمانی، برای استفاده از نظریه ایستایی لازم است که سری نایستا را به ایستا تبدیل کنیم. مثلاً می‌توانیم روند و تغییرات فصلی را از مجموعه داده‌ها حذف کرده و سپس به وسیله یک فرآیند تصادفی ایستا، تغییر در باقیمانده‌ها را الگوسازی کنیم. [۲]

الف) متوسط تعداد تقاضن روزانه پیدا شده در هر کامیون

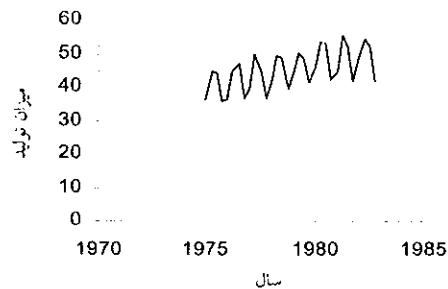


ج) تولید سه ماهه بستنی آمریکا

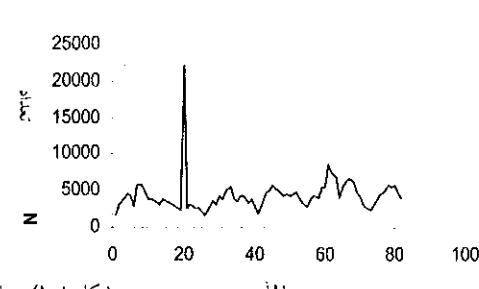
ب) تولید سالانه توتوون ایالات متحده



د) داده‌های مربوط به مگس گوشت



شکل ۱-۸(۱) چهار سری زمانی



شکل (۱-۸) که چهار سری زمانی را نشان می‌دهد، خصوصیات اجزاء سری‌های زمانی را آشکار می‌کند. به نظر می‌رسد که متوسط تعداد تقاضن روزانه پیدا شده در هر کامیون، در انتهای خط تولید کارخانه تولید کامیون که در شکل (۱-۸-الف) نشان داده شده است، در حول سطح ثابتی، نوسان می‌کند. همان‌طور که گفته شد سریهای زمانی که این پدیده را نشان می‌دهند، ایستاده در میانگین نامیده شده و حالتی ویژه سریهای زمانی ایستاده هستند.

تولید سالانه توتوون ایالات متحده که در شکل (۱-۸-ب) نشان داده شده در حول سطح ثابتی تغییر نمی‌کند، بلکه در کل یک روند رو به بالا را نشان می‌دهد. علاوه بر این، واریانس این سری توتوون، با اضافه شدن سطح سری، افزایش می‌یابد. سریهای زمانی که این پدیده را نشان می‌دهند، نایستا در میانگین و واریانس گفته شده و مثالهایی از سریهای زمانی نایستا هستند. تولید سه ماهه بستنی آمریکا در شکل (۱-۸-ج) طرح خاص دیگری را نشان می‌دهد که به واسطه تغییرات فصلی، طبیعتی تکراری دارد. سریهای زمانی که تغییرات فصلی را دربرمی‌گیرند، سریهای زمانی فصلی نامیده می‌شوند.

سریهای زمانی نایستا را مانند آنها بی که در شکل‌های (۸-۱-ب) و (۸-۱-ج) نشان داده شده‌اند، می‌توان با تبدیلات مناسبی به سری ایستا تبدیل نمود. سری زمانی چهارم که در شکل (۸-۱-د) نشان داده شده است، داده‌های مربوط به مگس گوشت است. این سری پدیده دیگری از نایستایی به‌دلیل تغییر ساختاری ناشی از یک یا چند اغتشاش خارجی را منعکس می‌کند. این نوع نایستایی را نمی‌توان با یک تبدیل استاندارد حذف نمود.

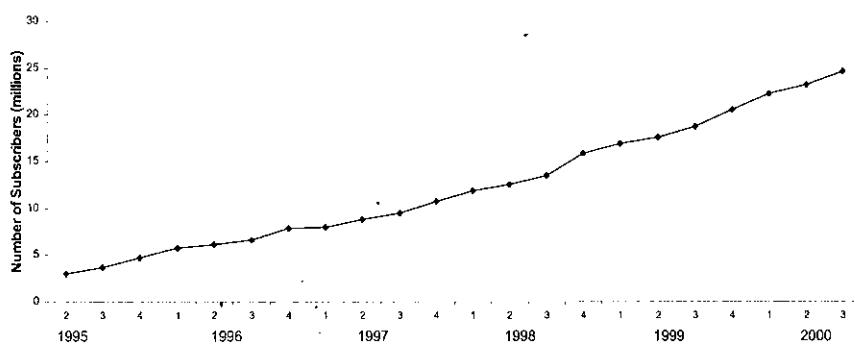
### ۸-۱-۲- شناسایی، تجزیه و حذف اجزاء سریهای زمانی

برای تجزیه و تحلیل مجموعه‌ای از داده‌ها، در اولین مرحله لازم است که نمودار مشاهدات را نسبت به زمان رسم کنیم. این کار غالباً مهم‌ترین خواص یک سری زمانی مانند روند، فصلی بودن و مشاهدات دورافتاده را آشکار می‌کند.

روشهای متعددی برای شناسایی، تعیین و یا حذف برخی از اجزاء سری‌های زمانی وجود دارد. با توجه به اینکه مهم‌ترین اجزاء یک سری زمانی، جزء روند و فصلی می‌باشد، در ادامه به بررسی روشهای شناسایی، هموارسازی، تجزیه و یا حذف این اجزاء پرداخته می‌شود. [۲]

### ۸-۱-۳- سریهای زمانی با روند خطی

اگر جزء روند در یک سری زمانی یک حالت افزایشی یا کاهشی مستقیم داشته باشد، می‌توان معادله روند را به شکل  $y_t = a + bt + e_t$  نشان داد که  $e_t$  جزء تصادفی در این رابطه است. شکل (۲-۸) معادله یک روند خطی را نشان می‌دهد.



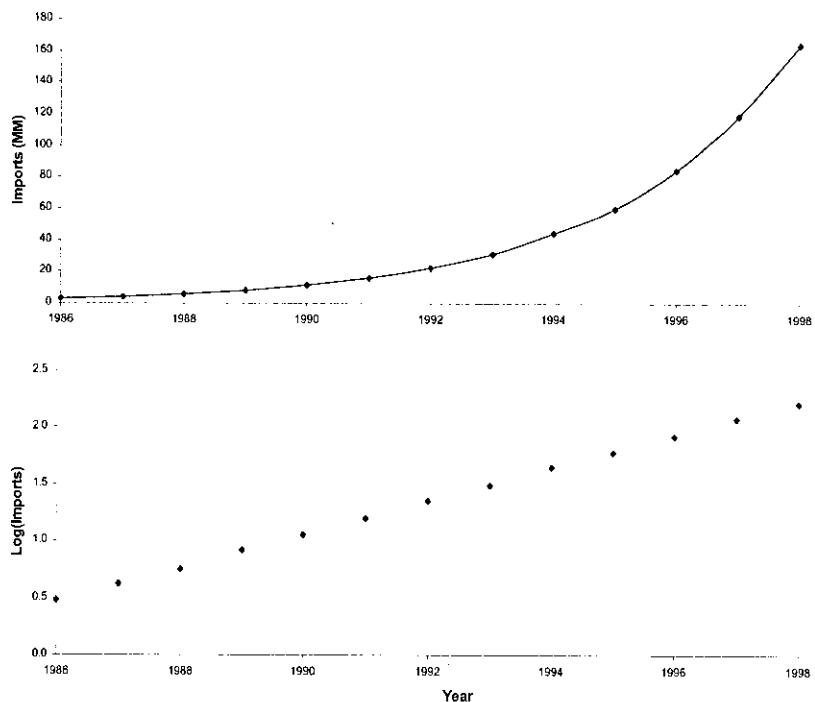
شکل (۲-۸) سری زمانی با روند خطی

### سریهای زمانی با روند غیرخطی

گاهی اوقات می‌توان با در نظر گرفتن تبدیلات خاصی، مانند لگاریتم یا ریشه دوم داده‌ها، روند غیرخطی سری‌ها را به یک روند خطی تبدیل کرد. در صورتی که انحراف معیار با میانگین نسبت مستقیم داشته باشد، یک تبدیل لگاریتمی مطابق روابط (۳-۸) و (۴-۸) مناسب است.

$$\log(y_t) = a + b + e_t \quad (3-8)$$

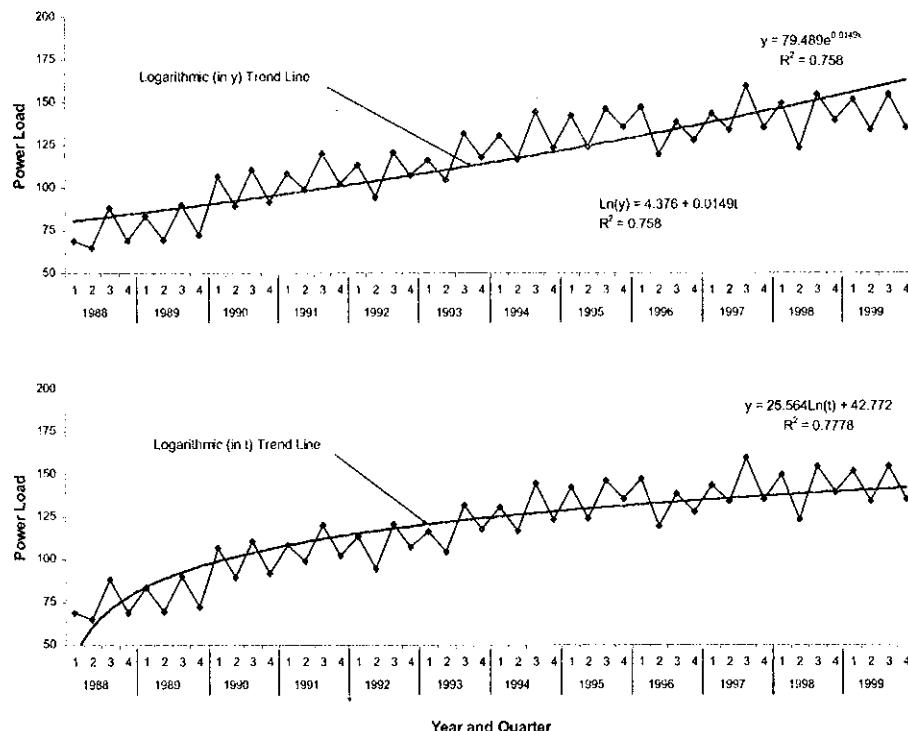
$$y_t = \exp(a + bt + e_t) \quad (4-8)$$



به ترتیب با روند غیرخطی و تبدیل یافته آن شکل (۳-۸) سری زمانی

اگر یک سری زمانی با یک نرخ کاهشی بر حسب زمان ظاهر شود، ممکن است رابطه (۵-۸) مناسب باشد.

$$y_t = a + b \ln(t) + e_t \quad (5-8)$$



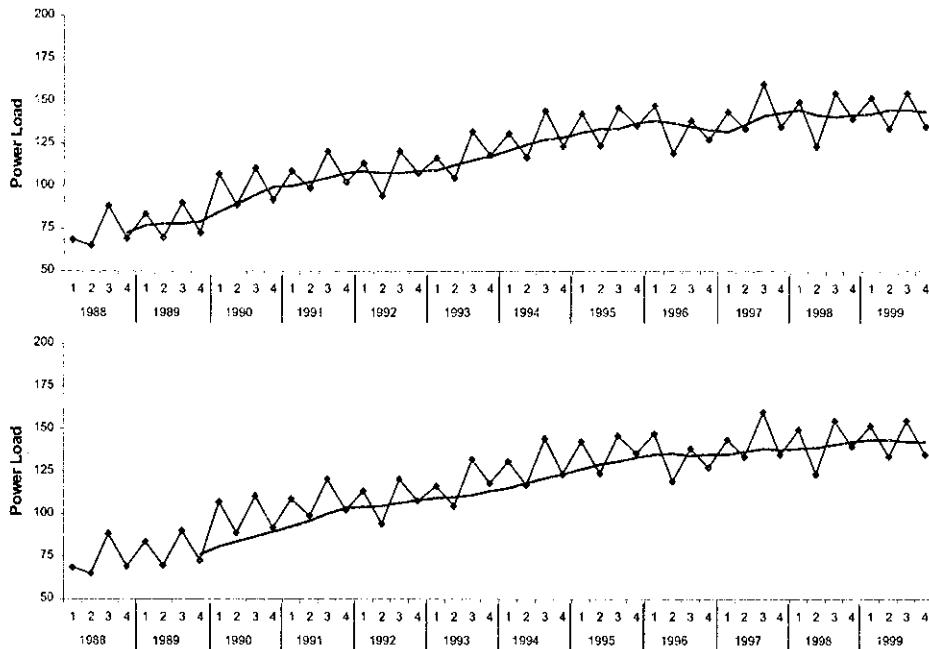
شکل ۸-۴) سری زمانی با تبدیل لگاریتمی روی متغیر وابسته (بالا) و زمان (پایین)

### میانگین متحرک

یک روش دیگر برای ارزیابی روند در سری‌های زمانی محاسبه میانگین  $m$  مشاهده اخیر می‌باشد. این روش به میانگین متحرک موسوم است و مطابق رابطه (۶-۸) محاسبه می‌گردد.

$$\bar{y}_{ma(t)} = \frac{(y_t + y_{t-1} + y_{t-2} + y_{t-3})}{4} \quad (6-8)$$

میانگین متحرک عمدها نوسانات موجود در داده‌ها را همواره می‌کند. این روش معمولاً زمانی به خوبی عمل می‌کند که داده‌ها یک روند خطی و الگوی منظمی از نوسانات داشته باشند.



شکل ۵-۸) روش میانگین متحرک در سری زمانی به ترتیب چهار نقطه و هشت نقطه

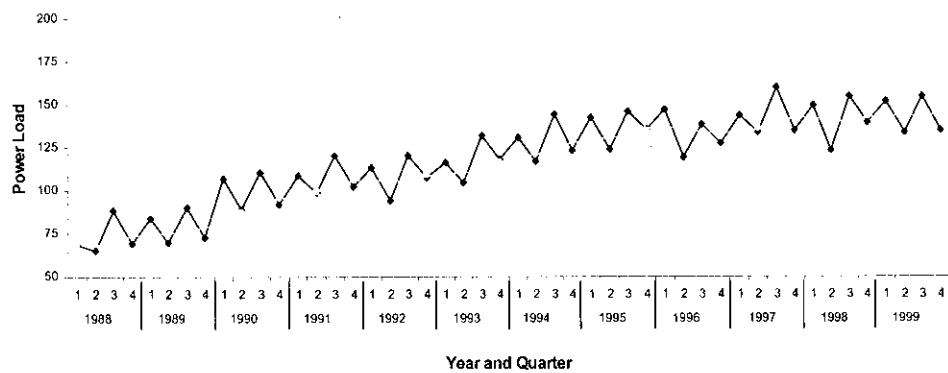
### هموارسازی نمایی

مهم‌ترین عیب روش میانگین متحرک این است که وزن یا اهمیت داده‌های گذشته، به صورت یکسان در نظر گرفته می‌شود. برای برطرف کردن این اشکال از یک روش جامع‌تر که یک روش میانگین متحرک موزون است و در آن تخصیص وزن به دوره‌های گذشته به صورت یک تصاعد هندسی می‌باشد، استفاده می‌کنیم. این روش هموارسازی نمایی نام دارد و به مقادیر گذشته سری زمانی تا به حال، وزن داده می‌شود و البته وزن بیشتری برای داده‌های جدید در نظر گرفته شده و هرچه به سمت داده‌های قدیمی پیش می‌رویم، وزن آنها طی یک فرآیند نمایی، کاهش می‌یابد.

<p>Let <math>w=0.5</math></p> $S_1 = Y_1$ $S_2 = 0.5Y_2 + (1-0.5)S_1 = 0.5Y_2 + 0.5Y_1$ $S_3 = 0.5Y_3 + (1-0.5)S_2 = 0.5Y_3 + 0.25Y_2 + 0.25Y_1$ $S_4 = 0.5Y_4 + (1-0.5)S_3 = 0.5Y_4 + 0.25Y_3 + 0.125Y_2 + 0.125Y_1$	$S_1 = Y_1$ $S_t = wY_t + (1-w)S_{t-1}$ $= wY_t + w(1-w)Y_{t-1} + w(1-w)^2Y_{t-2} + \dots$
---	--

شکل ۶-۸) هموارسازی میانگین متحرک (راست) و میانگین متحرک موزون نمایی (چپ)

هرچه مقدار ثابت هموارسازی ( $w$ ) بیشتر باشد نشان دهنده آن است که داده‌های قدیمی اثر کمتری بر روی هموارسازی پیش‌بینی دارند و هرچه مقدار  $w$  کمتر باشد داده‌های سری زمانی هموارتر خواهند بود.

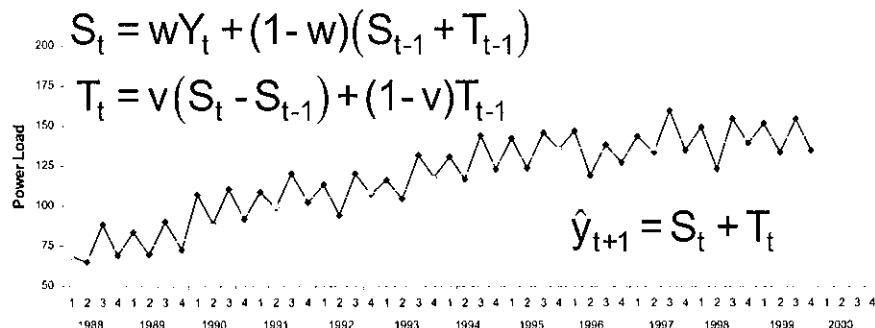


شکل ۸-۸) میانگین متحرک موزون نمایی ( $w=0.34$ )

انتخاب مقدار  $w$  را می‌توان براساس حداقل کردن شاخصهای میانگین قدر مطلق خطأ و میانگین خطای مریع شده که دربحث پیش‌بینی به آنها اشاره شد، انجام داد. در مثال بالا  $W=0.34$ ، براساس حداقل کردن مقدار میانگین خطای مریع شده به دست آمده است.

### هموارسازی نمایی با تنظیم روند

مانند روش میانگین متحرک، هموارسازی نمایی ساده، نسبت به روند واکنش نشان نمی‌دهد. در این حالت استفاده از روش‌های دیگر نظری هموارسازی نمایی با تنظیم روند مناسب است.



شکل ۸-۸) پیش‌بینی با استفاده از میانگین متحرک موزون نمایی ( $w=0.34$ )

### هموارسازی نمایی با تنظیم اثر روند و فصلی

همان‌طور که گفته شد، در تجزیه و تحلیل سریهای زمانی گاهی اوقات لازم است اثرات روند فصلی یا سیکلی را به‌طور مجزا مطالعه کنیم. برای حذف کردن و یا اندازه‌گیری هریک از این مؤلفه‌ها، در تحلیل سریهای زمانی از روش‌های مربوطه استفاده می‌شود. روش‌های میانگین متحرک و غیره که قبلاً به آنها اشاره شد، برخی از این روشها می‌باشند. یکی دیگر از روش‌های مورد استفاده برای حذف جزء روند یا فصلی، استفاده از روش تفاضلی است. مثلاً برای داده‌های غیر فصلی، تفاضل مرتبه اول طبق رابطه (۷-۸) برای رسیدن به ایستایی کافی است.

$$y_t = x_{t+1} - x_t \quad (7-8)$$

همچنین برای برآورد اثر فصلی برحسب اینکه اثر فصلی جمعی یا ضربی باشد، می‌توان با استفاده از روابط (۸-۸) تا (۹-۸) آن را برآورد کرد.

$$x_t - S_m(x_t) \quad (8-8)$$

$$\frac{x_t}{S_m(x_t)} \quad (9-8)$$

$$S_m(x_t) = \frac{\frac{1}{2}x_{t-1} + x_{t-2} + \dots + x_{t-5} + \frac{1}{2}x_{t+1}}{12} \quad (10-8)$$

همچنین یکی از راههای حذف کردن اثر فصلی، استفاده از شیوه تفاضلی است. [۲]

#### ۴-۱-۸- جستجوی تشابه در تحلیل سریهای زمانی

مطابق آنچه که در ابتدای فصل اشاره شد، مباحثت مدلسازی و پیش‌بینی سریهای زمانی در بسیاری از مراجع آماری مورد بررسی قرار گرفته است. لیکن آماردانها روش‌های مناسبی را برای تشابه و شاخص‌گذاری سریهای زمانی بررسی نکرده‌اند. تعداد زیادی از این مسائل توسط جامعه علمی کامپیوتر حل شده‌اند. یکی از مسائل جالب در داده‌های سری زمانی، یافتن سریهای زمانی متفاوتی است که رفتار مشابه داشته باشند. مسئله را می‌توان این‌گونه نیز مطرح کرد که آیا دو سری زمانی  $X$  و  $Y$  داده شده مشابه هستند یا خیر؟ به عبارت دیگر، تابع  $Sim(X, Y)$  را تعریف کرده و میزان تشابه دو سری یا به‌طور مشابه، تابع فاصله  $Dis(X, Y)$  را محاسبه می‌کنیم. برای نمونه، هر سری زمانی مقدار و سیر تدریجی یک شیء را به صورت تابعی از زمان در مجموعه‌ای از داده‌های جمع‌آوری شده توصیف می‌کند (مثلاً قیمت سهام). هدف می‌تواند خوشه‌بندی اشیاء مختلف در گروه‌های مشابه (مانند گروهی از سهامها که تغییر قیمت یکسان داشته‌اند) یا دسته‌بندی اشیاء براساس مجموعه‌ای از ویژگیهای شناخته شده باشد. این مورد مشکل است، چون مدل تشابه باید اجازه تطابق غیر دقیق را بدهد.

یکی از مسائل جالب مطرح شده در تشابه دنباله‌ها<sup>۱</sup>، تشابه زیر دنباله‌ها می‌باشد که در آن برای یک سری زمانی داده شده  $X$  و الگوی سری زمانی کوتاه‌تر  $Y$ ، می‌خواهیم زیر دنباله‌ای از  $X$  را که مشابه الگوی  $Y$  عمل می‌کند پیدا کنیم. برای پاسخ به این پرسش، نظریه‌های متفاوتی از تشابه سریهای زمانی در پژوهش‌های داده‌کاوی مطرح شده است. در این بخش مدل‌های مختلف اندازه‌گیری تشابه سریهای زمانی مطرح می‌شود که براساس شاخصهای کارآیی و دقت، می‌توان آنها را ارزیابی کرد. نمونه‌هایی از اندازه‌گیری تشابه را که بر اساس نرم اقلیدسی، تخمینهای خطی قطعه‌ای<sup>۲</sup>، تاباندن زمانی پویا<sup>۳</sup> (DTW) و بزرگترین زیر دنباله‌های مشترک<sup>۴</sup> (LCSS) هستند را بررسی خواهیم کرد. [۵]

<sup>۱</sup>- Sequences

<sup>۲</sup>- Piecewise Linear Approximations: PLA

<sup>۳</sup>- Dynamic Time Warping: DTW

<sup>۴</sup>- Longest Common Subsequences Similarity: LCSS

## ۱-۵- مقیاسهای اندازه‌گیری تشابه در سربهای زمانی

### فاصله‌اقلیدسی و نرم<sup>۱</sup>

یکی از ساده‌ترین راههای اندازه‌گیری تشابه در سربهای زمانی اندازه‌گیری فاصله اقلیدسی است. دو دنباله زمانی با طول  $n$  را فرض کنید. ما هر دنباله را در فضای  $n$  بعدی اقلیدسی، به عنوان یک نقطه می‌بینیم. عدم شباهت یا فاصله بین دو دنباله  $X$  و  $Y$  را با  $L_p(X, Y)$  تعریف می‌کنیم (وقتی  $p=2$  است، این فاصله، همان فاصله اقلیدسی معروف است). این اندازه‌گیری مزایای مختلفی دارد. فهم آن آسان، محاسبه آن ساده و برای حل مشکلات دیگر مثل شاخص‌گذاری و خوشبندی سربهای زمانی قابل استفاده است. هرچند معایب زیادی نیز دارد که آن را برای کاربردهای متعددی نامناسب می‌کند. یکی از اصلی‌ترین معایب آن این است که اجازه نمی‌دهد که دنباله‌های زمانی خط مبنای متفاوتی داشته باشند. برای مثال سهام  $X$  با نوسان حدود \$100 و  $Y$  با نوسان حدود \$30 را در نظر بگیرید. حتی اگر شکل هر دو دنباله زمانی به هم خیلی شبیه باشد، ممکن است فاصله اقلیدسی بین آنها خیلی زیاد شود. همچنین نمی‌توان به این روش در مقیاسهای مختلف اندازه‌گیری نمود. برای مثال ممکن است سهام  $X$  در دامنه کوچکی نوسان کند (بین \$95 تا \$105) درحالی که سهام  $Y$  در دامنه بزرگتری نوسان کند (بین \$20 تا \$40).

### تبديلات نرمال<sup>۲</sup>

با استفاده از نرمال کردن دنباله‌ها می‌توان معایب نرم  $L_p$  را در اندازه‌گیری تشابه برطرف کرد. در رابطه (۱۰-۸) اگر  $(X, \mu)$  میانگین و  $(\sigma, \sigma)$  واریانس دنباله باشد، دنباله  $\{x_1, \dots, x_n\} = X$  را با دنباله نرمال  $X'$  جایگزین می‌کنیم:

$$x'_i = (x_i - \mu(X)) / \sigma(X) \quad (11-8)$$

همچنین دنباله  $Y$  را با دنباله نرمال  $Y'$  جایگزین می‌کنیم. در نهایت عدم تشابه بین  $X$  و  $Y$  را با  $L_p(X', Y')$  تعریف می‌کنیم. این تعریف تشابه، معایب استفاده مستقیم از نرم  $L_p$  را که دنباله‌ها نرمال نیستند، حل می‌کند. برای مثال دو سهام  $X$  و  $Y$  که قبلاً در مورد آن صحبت شد،

<sup>۱</sup>- Euclidean Distances and Lp Norms

<sup>۲</sup>- Normalization Transformations

را در نظر بگیرید. بعد از نرمال شدن هر دو دارای یک خط مبنا شده (چون میانگین هر دو با نرمال کردن یکسان می‌شود) و دامنه پکسانی خواهد داشت (چون با نرمال کردن واریانس داده‌ها یکسان شده است).

فرآیند نرمال کردن هم معایب خاص خود را دارد. مثلاً به تأخیر و تقدم فاز در زمان خیلی حساس است. به عنوان نمونه دو دنباله  $X$  و  $Y$  را در نظر بگیرید که  $X$  شبیه به موج سینوسی است، در صورتی که  $Y$  شبیه به موج کسینوسی است. هر دو اصولاً دارای شکل یکسانی هستند به جز اینکه یک تأخیر فاز دارند. ولی فاصله بین این دو دنباله قابل ملاحظه است (در هردو حالت نرمال شده یا نرمال نشده). همچنین فاصله‌اقلیدسی، شتاب و کاهش شتاب در طول محور اصلی را نشان نمی‌دهد. برای مثال دو دنباله  $X$  و  $Y$  که شبیه موج سینوسی هستند را در نظر بگیرید که فقط پریود  $X$  دو برابر دوره  $Y$  است. حتی اگر این دو دنباله نرمال شوند، فاصله‌اقلیدسی نمی‌تواند شباهت بین این دو سیگنال را نشان دهد [۵].

### تبديلات عمومی<sup>۱</sup>

تشخيص اهمیت نظریه حالت، در محاسبات تشابه در سال ۱۹۹۵ مطرح شد. یک چارچوب تشابه عمومی شامل زبان قواعد تبدیلات توصیف شده است. هر روش در زبان تبدیلات، یک دنباله ورودی را گرفته و با هزینه مربوط به آن روش، دنباله خروجی متناظرش را تولید می‌کند. تشابه بین دنباله  $X$  و  $Y$ ، حداقل هزینه برای تغییر  $X$  به  $Y$  با استفاده از چنین روش‌هایی می‌باشد. برای مثال در شکل زیر دنباله‌هایی که به شکل منحنی‌های خطی-قطعه‌وار هستند، نشان داده شده است. به عنوان نمونه، یک قاعده تبدیل می‌تواند یکی کردن بخش‌های مجاور روی یک بخش باشد. هزینه این قاعده می‌تواند تابعی از طول و شبیه بخش جدید و بخش ابتدایی باشد. قاعده دیگر می‌تواند یک بخش منفرد را با یک جفت بخش مجاور جایگزین کند.

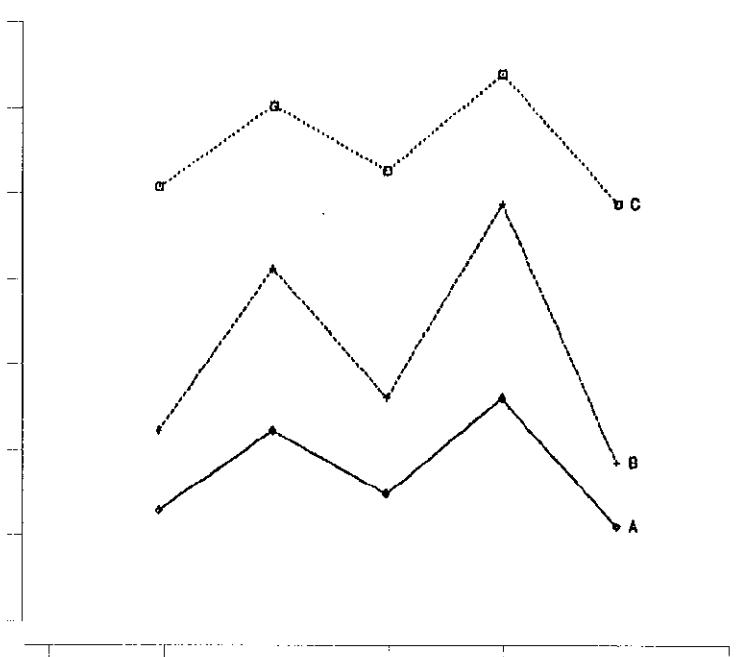


شکل ۹-۸) دنباله‌هایی به شکل منحنی‌های خطی

یکی از قواعد تبدیلات استفاده از میانگین متحرک برای هموارکردن سریهای زمانی است. برای مثال سری زمانی  $X$  که قیمت روزانه سهام در طول یک فاصله زمانی چند روزه است را در نظر می‌گیریم. برای هر ۳ روز یک میانگین متحرک حساب کرده و در دنباله  $X'$  تغییر می‌دهیم و نتیجه دنباله  $X'$  می‌شود که  $x' = (x_{i-1} + x_i + x_{i+1})/3$ .

یکی دیگر از روش‌های تبدیل، تبدیل مقیاس و انتقال است. در تبدیل مقیاس، هر جزء با مقیاس ثابتی افزایش می‌یابد. برای مثال هر  $x_i$  با  $cx_i$  جایگزین می‌شود که  $c$  مقداری ثابت است. هر تبدیل انتقال هر جزء را با یک عدد ثابت از موقعیت فعلی به سمت راست یا چپ منتقل می‌کند (هر  $x_i$  با  $x_i+c$  جایگزین می‌شود که  $c$  یک عدد صحیح ثابت است).

مثال: سه دنباله شکل زیر را در نظر بگیرید.



شکل ۱۰-۸) سه دنباله نمونه

$$A = (5, 10, 6, 12, 4)$$

$$B = (10, 20, 12, 24, 8)$$

$$C = (25, 30, 26, 32, 24)$$

این سه سری زمانی متفاوتند، اما باهم رابطه نزدیکی دارند. سری A می‌تواند با دو برابر کردن جملات به B تبدیل شود و C با ۲۰ واحد انتقال می‌تواند به A تبدیل شود. به علاوه B می‌تواند با نصف کردن جملاتش و سپس ۲۰ واحد انتقال تبدیل به C شود. این بدین معنی است که این سریها با تبدیل مقیاس‌بندی و انتقال مناسب، در واقع یکی هستند. اگر سه دنباله بالا را به عنوان روند قیمت سه سهام در نظر بگیریم، با وجود اینکه قیمت سهام شرکت C بیشتر از شرکت A است، اما چون نوسان یکسانی دارند دقیقاً از روند قیمتی یکسانی پیروی می‌کنند. یا اگر چه قیمت سهام شرکت B همیشه دو برابر قیمت سهام شرکت A ولی نوسان آنها متناسب با قیمتیان است در نتیجه روند قیمتی آنها باید یکسان در نظر گرفته شود. سری A با سری B مشابه است، اگر A بتواند با یکی از تبدیلات گفته شده به B تبدیل شود [۷].

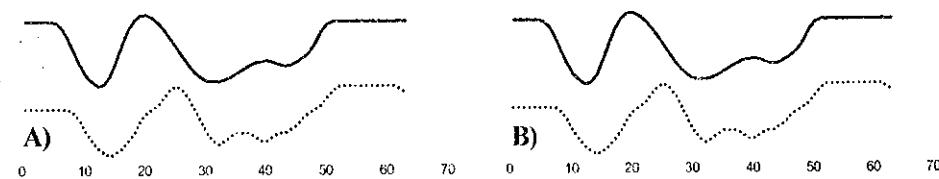
قواعد تبدیل، یک روش عمومی را برای تعریف تشابه که مناسب کاربردهای خاصی است، پیشنهاد می‌دهد. هرچند بعضی از معاایب آنها، مشکلاتی نیز ایجاد می‌کنند. مثلًا محاسبات زیر دنباله‌ها (مانند شاخص‌گذاری) مشکل می‌شود، چون استخراج مشخصه‌ها از دنباله X، مخصوصاً اگر قواعد استفاده شده به هر دو دنباله X و Y بستگی داشته باشد، کار پیچیده‌ای است. همچنین فاصله‌اقلیدسی در فضای مشخصه ممکن است تقریب خوبی برای عدم تشابه دنباله‌های ابتدایی نباشد [۵].

## ۱-۸- تاباندن محور زمان به صورت پویا<sup>۱</sup>

یکی از معمول‌ترین کارهایی که با داده‌های سریهای زمانی انجام می‌دهند، مقایسه یک دنباله با دنباله دیگر است. در بعضی از کاربردها یک مقیاس اندازه‌گیری ساده مانند اندازه‌گیری فاصله‌اقلیدسی، کافی است. اگرچه اغلب این حالت پیش می‌آید که دو دنباله با اجزای تقریباً یکسان در محور Xها، در قسمتی نسبت به هم کشیده‌تر هستند. شکل زیر این موضوع را با یک مثال ساده نشان می‌دهد. برای فهمیدن شباهت چنین دنباله‌هایی، قبل از میانگین گرفتن از آنها،

<sup>۱</sup>- Dynamic Time Warping

به عنوان یک قدم پیش‌پردازش، یک یا هر دو دنباله را روی محور زمان می‌پیچانیم. روش  $DTW$  برای این نوع تاباندن روی محور زمان، کارآمد می‌باشد. علاوه بر داده‌کاوی،  $DTW$  در تشخیص حرکات علم روباتیک، تحلیل ساختاری و پزشکی کاربرد دارد.



شکل (۱۱-۸) دو دنباله که وضعیت دستخط یک شخص را هنگام نوشتن کلمه Pen در زبان علامت روی محورها نشان می‌دهد.

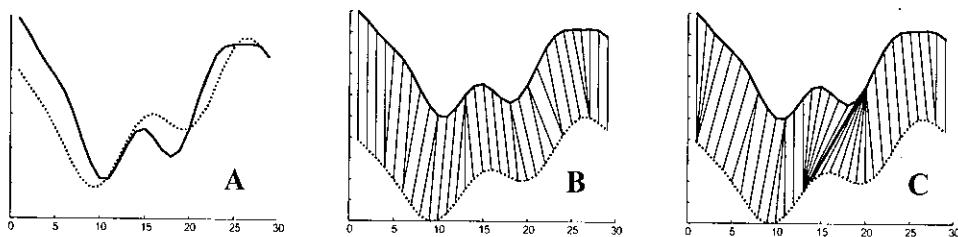
شکل (۱۱-۸) یک مثال از کاربرد  $DTW$  می‌باشد. نمودار A شامل دو دنباله است که وضعیت دستخط یک شخص را هنگام نوشتن کلمه Pen در زبان علامت روی محورها نشان می‌دهد. دنباله‌ها در دو روز مختلف ضبط شده‌اند. توجه کنید با وجود اینکه دنباله‌ها شکل عمومی یکسانی دارند، ولی در محور زمان بر هم منطبق نیستند. یک مقیاس فاصله که فرض می‌کند نامین نقطه روی یک دنباله منطبق بر نامین نقطه روی دنباله دیگر است، باعث یک عدم تشابه نامید کننده می‌شود. نمودار B می‌تواند به صورت کارآمد یک تطابق بین دو دنباله ایجاد کند که محاسبه فاصله فاصله پیچیده‌تری را ایجاد می‌کند.

اگر چه استفاده از  $DTW$  در بسیاری از زمینه‌ها موفق بوده است، می‌تواند نتایج غیر قابل کنترل و نامطلوبی داشته باشد. مشاهدات حاکی از آن است که الگوریتم می‌تواند تغییرات در محور  $x$  را با تاباندن روی محور  $x$  ها توضیح دهد. این می‌تواند باعث تطابق‌های غیرشهودی هنگام تصویر یک نقطه منفرد از یک سری زمانی، روی یک زیربخش بزرگ از سری زمانی دیگر شود. ما چنین رفتار غیردلخواهی را مقادیر منفرد یا تکینهای<sup>۱</sup> می‌نامیم. بخش وسیعی از روشها برای مقابله با این رفتارها ارائه شده‌اند. دستاورد این روشها، چگونگی تاباندن مجاز را مشخص

می‌کنند. اگر چه استفاده از این روشها در بعضی مواقع از یافتن روش تاباندن صحیح جلوگیری می‌کند.

در موارد شبیه‌سازی شده، تاباندن وقتی صحیح تشخیص داده می‌شود که ما ابتدا یک سری زمانی را بتابانیم و سپس سعی کنیم سری اصلی را از روی سری تابانده شده به دست آوریم. در رویدادهای طبیعی، منظور ما از روش تاباندن صحیح آن است که مانند شکل (۱۲-۸-B) به طور شهودی تطابق یک به یک مشخصه‌ها واضح باشد.

یک مشکل دیگر با DTW آن است که الگوریتم ممکن است در پیدا کردن تطابق‌های طبیعی و واضح در دو دنباله، تنها به این دلیل که یک مشخصه در یک دنباله کمی بالاتر یا پایین‌تر از مشخصه مربوطه در دنباله دیگر است، اشتباہ کند. برای مثال نقاط اوج یا حضیض، نقطه خمیدگی، قسمت مسطح و غیره شکل (۱۲-۸) این مسئله را نشان می‌دهد.



شکل ۱۲-۸- (A) دو سیگنال ترکیبی (با میانگین و واریانس یکسان). (B) تطابق نظیر به نظری مشخصه‌ها. (C) تطابق ایجاد شده به وسیله DTW

توجه کنید که DTW، دو نقطه اوج مرکزی را به علت اینکه آنها در محور Y‌ها کمی با هم فاصله دارند، منطبق در نظر گرفته است [۸].

### ۷-۱-۸- الگوریتم کلاسیک DTW

فرض کنید دو سری زمانی Q و C را با طول‌های m و n داریم [۸] :

$$Q = q_1, q_2, \dots, q_i, \dots, q_n$$

$$C = c_1, c_2, \dots, c_j, \dots, c_m$$

برای تطابق دو دنباله با استفاده از DTW یک ماتریس  $m \times n$  می‌سازیم به طوری که عنصر  $(ith)$  ماتریس، شامل فاصله  $d(q_i, c_j)$  بین دو نقطه  $c_j$  و  $q_i$  می‌باشد (معمولًاً از فاصله اقلیدسی

استفاده می‌کنیم)، بنابراین  $d(q_i, c_j) = (q_i - c_j)^T d$ . هر عنصر  $(i, j)$  ماتریس به تطابق بین نقاط  $q_i$  و  $c_j$  مربوط است که در شکل (۱۲-۸) نشان داده شده است. یک مسیر تاباندن  $W$ ، یک دنباله از عناصر پیوسته ماتریس است که یک نگاشت را بین  $Q$  و  $C$  مشخص می‌کند. عنصر  $k$  ام  $W$  به صورت  $w_k = (i, j)$  تعریف می‌شود و در نتیجه خواهیم داشت.

$$W = w_1, w_2, \dots, w_k, \dots, w_n \quad \max(m, n) \leq K < m + n - 1 \quad (12-8)$$

مسیر تاباندن معمولاً محدودیتهایی دارد.

**شرایط حدی:**  $w_1 = (m, n)$  و  $w_k = (1, 1)$  به راحتی نشان می‌دهند که مسیر تاباندن از اولین نقطه روی قطر اصلی شروع و به نقطه مقابل آن در انتهای قطر اصلی خاتمه می‌یابد.

**پیوستگی:** اگر  $w_k = (a, b)$  و  $w_{k-1} = (a', b')$  باشند، لازم است که  $a - a' \leq 1$  و  $b - b' \leq 1$ . این شروط گامهای مجاز را در مسیر تاباندن برای سلوهای مجاور از جمله سلوهای مجاور قطری مشخص می‌کند (هیچ جزئی نمی‌تواند در دنباله حذف شود).

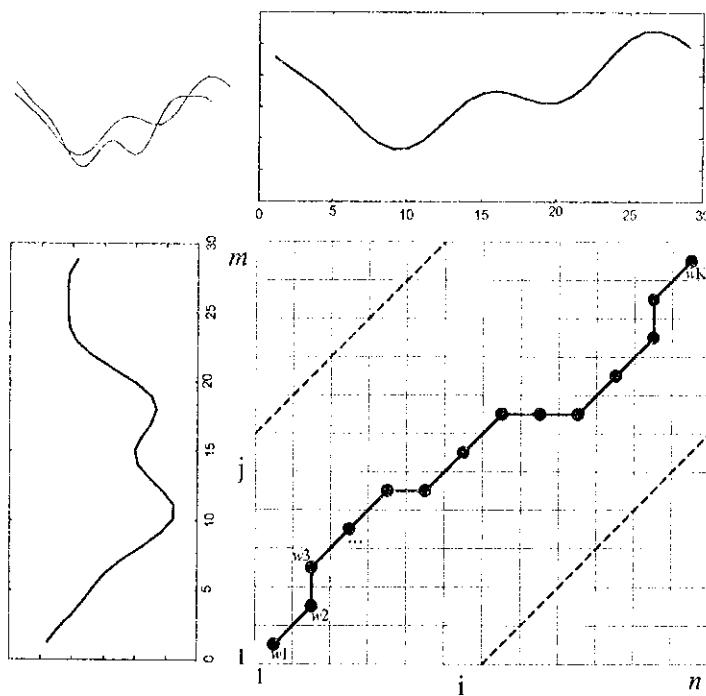
**یکنواختی:** اگر  $w_k = (a, b)$  و  $w_{k-1} = (a', b')$  لازم است که  $a - a' \geq 0$  و  $b - b' \geq 0$ . این باعث می‌شود نقاط در  $W$  به صورت یکنواخت و بر حسب زمان قرار بگیرند.

باتوجه به محدودیتهای گفته شده تعداد زیادی مسیر تاباندن وجود دارد ولی تنها مسیرهایی که هزینه تاباندن را حداقل می‌کنند مورد نظر ما هستند.

$$DTW(Q, C) = \min \left\{ \sqrt{\sum_{k=1}^K w_k} / K \right\} \quad (13-8)$$

برای نشان دادن اینکه مسیرهای مختلف می‌توانند طول‌های مختلف داشته باشند،  $K$  در رابطه (۱۲-۸) در مخرج ظاهر می‌شود. این مسیر را می‌توان به صورت کارآمد با استفاده از برنامه‌ریزی پویا به دست آورد. برای این کار  $\gamma(i, j)$  را که فاصله تجمعی است به دست می‌آوریم.  $d(i, j)$  در فرمول (۱۳-۸) همان فاصله  $i$  و  $j$  است که در ماتریس به دست آورده‌یم و قسمت بعد حداقل مقدار فاصله سلوهای مجاور است.

$$\gamma(i, j) = d(q_i, c_j) + \min \{ \gamma(i-1, j-1), \gamma(i-1, j), \gamma(i, j-1) \} \quad (14-8)$$



شکل ۱۳-۸) مثالی از یک مسیر تاباندن

### محدودیت‌های الگوریتم کلاسیک DTW

مشکل مقادیر تکین از اوایل سال ۱۹۷۸ توسط ساکوئی و چیبا<sup>۱</sup> مورد توجه قرار گرفت. روش‌های مختلفی برای کم رنگ‌تر کردن این مشکل مطرح شده است که ما اجمالاً به بررسی آنها می‌پردازیم [۸].

۱) پنجره بندی<sup>۲</sup>: عناصر مجاز ماتریس به آنها بین که در پنجره می‌افتد، محدود می‌شوند. مشکل مقداری تکین از  $|i - (n/(m/j))| < R$  یک عدد صحیح مثبت است و نمایانگر عرض پنجره است. این بدین معنی است که گوشش‌های ماتریس هرس می‌شوند. همان‌طور که در شکل (۱۳-۸) با خط‌چین نشان داده شده است. دیگران پنجره‌بندی را با اشکال مختلف دیگری تجربه کردند این دستاورد تا جای ممکن مشکل تکینها را محدود می‌کند اما از رخداد آن جلوگیری نمی‌کند.

<sup>۱</sup>- Sakoe & Chiba

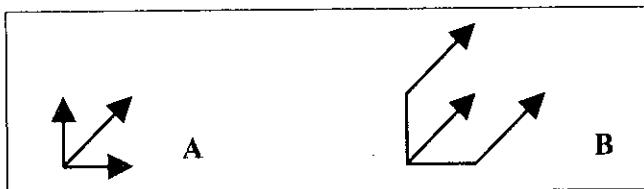
<sup>۲</sup>- Windowing

۲) وزن دهنی به شیب<sup>۱</sup>: اگر در معادله (۱۵-۸) فاصله را به صورت زیر در نظر بگیریم:

$$\gamma(i, j) = d(i, j) + \min\{\gamma(i-1, j-1), X\gamma(i-1, j), X\gamma(i, j-1)\} \quad (15-8)$$

به طوری که  $X$  یک عدد حقیقی مثبت است و می توانیم شکل تاب دادن را با عوض کردن مقدار  $X$  محدود کنیم. هرچه مقدار  $X$  بزرگتر شود مسیر تاب دادن به سمت قطعی شدن تمایل پیدا می کند.

۳) محدودیت های شیب<sup>۲</sup>: ما می توانیم معادله فاصله را به صورت یک نمودار الگوی گامهای قابل قبول به تصویر درآوریم. مثلاً در شکل (۴-۱۴-۸) پیکانها نشان دهنده گامهای مجازی هستند که مسیر تاباندن در هر مرحله می تواند بردار را نشان دهد. می توانستیم معادله (۱۵-۸) را با معادله های زیر که به الگوی گام نشان داده شده در شکل (B-۱۴-۸) مربوط است جایگزین کنیم با استفاده از این معادله، مسیر تاباندن یک قدم قطعی، گام برمی دارد که این قدم می تواند بعد از طی یک قدم موازی با یکی از محورها صورت گیرد.



شکل ۱۴-۸) یک نمایش تصویری از رو نوع الگوی مسیر متفاوت

$$\gamma(i, j) = d(i, j) + \min(\gamma(\bar{i}-1, \bar{j}), \gamma(i-1, j), \gamma(i, j-1)) \quad \text{الگوی (A)}$$

$$\gamma(i, j) = d(i, j) + \min[\gamma(i-1, j-1), \gamma(i-1, j-2), \gamma(i-2, j-1)] \quad \text{الگوی (B)}$$

همه موارد گفته شده می توانند در محدود کردن مشکل تکینها کمک کنند ولی ریسک از دست دادن روش تاب دادن صحیح، هنوز وجود دارد. مشکل دیگر آن است که چگونگی انتخاب پارامترهای موجود در الگوها هنوز برای ما واضح نیست. برای مثال نمی دانیم عدد  $R$  را در پنجره بندی و عدد صحیح  $X$  را در الگوی وزن دهنی شیب چگونه به دست آوریم.

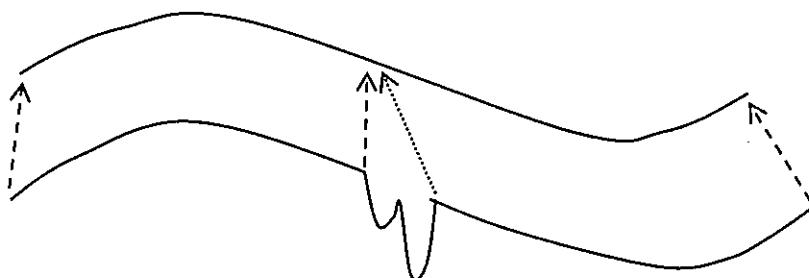
<sup>۱</sup>- Slope Weighting

<sup>۲</sup>- Slope Constraints

### ۸-۱-۸- شباهت بزرگترین زیردنباله مشترک (LCSS)

شباهت بزرگترین زیردنباله مشترک برای اندازه‌گیری اختلاف سریهای زمانی، در مواردی مانند تشخیص صدا و تطابق الگوی متن به کار می‌رود. ایده اصلی، تطبیق دو دنباله بر یکدیگر است، حتی اگر چند جزء آنها منطبق نباشد. روش LCSS دو مزیت دارد: (الف) بعضی از اجزا می‌توانند منطبق نباشد (مثل نقاط پرت)، ولی در فاصله اقلیدسی و DTW همه اجزاء حتی نقاط پرت باید منطبق باشد. (ب) همان‌طور که در ادامه دیده خواهد شد اندازه‌گیری LCSS، کارآیی محاسبات تقریبی را بیشتر می‌کند. حال به بررسی یک مثال شهودی‌تر می‌پردازیم. دو دنباله  $X$  و  $Y$  را در نظر بگیرید:

$$Y = \{2, 5, 4, 7, 3, 10, 8\} \quad , \quad X = \{3, 2, 5, 7, 4, 8, 10, 7\}$$



شکل (۱۵-۸) اندازه‌گیری LCSS

شکل (۱۵-۸) ایده اصلی اندازه‌گیری LCSS را مشخص می‌کند. فرض کنید قسمتهای مشخصی از هر دنباله در فرآیند تطابق حذف شده باشد (مثل نقاط پرت یا داده‌های مغشوش). مقدار LCSS بین  $X$  و  $Y$  برابر  $\{2, 5, 7, 10\}$  است.

دو دنباله  $X$  و  $Y$ ، به ترتیب با طول  $m$  و  $n$  را در نظر می‌گیریم. مشابه همان کاری که در DTW انجام شد، یک رابطه بازگشتی برای طول LCSS دنباله  $X$  و  $Y$  به دست می‌آید. مقدار  $L(i, j)$  نشان دهنده LCSS زیر دنباله  $\{x_1, \dots, x_i\}$  و  $\{y_1, \dots, y_j\}$  می‌باشد. ( $L(i, j)$  می‌تواند یک رابطه بازگشتی به صورت زیر باشد.

$$\text{If } x_i = y_j \text{ then } L(i, j) = 1 + L(i-1, j-1), \text{ else } L(i, j) = \max\{D(i-1, j), D(i, j-1)\}$$

ما عدم تشابه بین  $X$  و  $Y$  را با  $LCSS(X, Y) = (m+n-2l)/(m+n)$  تعریف می‌کنیم، که اطول  $LCSS$  است. این کمیت، حداقل مقدار نرمال شده تعداد عناصری است که باید از  $X$  حذف و یا به  $X$  اضافه شوند تا  $X$  به  $Y$  تبدیل شود. مشابه  $DTW$  اندازه‌گیری  $LCSS$  به وسیله برنامه‌ریزی پویا در زمان  $O(mn)$ ، می‌تواند محاسبه شود. اگر پنجره تطابق با طول  $w$  مشخص شده باشد که  $|i - j| \leq w$  می‌تواند حداکثر  $w$  باشد، می‌توان پیچیدگی آن را به مقدار  $O(wn)$  کاهش داد. الگوریتم یافتن طولانی‌ترین مسیر مشترک بین دو دنباله به صورت زیر است:

```

LCSS - LENGTH( $X, Y$ )
 $m = length[X], n = length[Y]$ 
 $for i = 1 to m do c[i,] \leftarrow .$ 
 $for j = 1 to n do c[,j] \leftarrow .$ 
 $for i = 1 to m$ 
 $do for j = 1 to n$ 
 $do if x_i = y_j$ 
 $then c[i, j] \leftarrow c[i - 1, j - 1] + 1 b[i, j] = "R"$ 
 $else if c[i - 1, j] \geq c[i, j - 1]$ 
 $then c[i, j] \leftarrow c[i - 1, j] b[i, j] = "↑"$ 
 $else c[i, j] \leftarrow c[i, j - 1] b[i, j] = "←"$ 

```

در  $LCSS$  نیاز به اینکه داده‌های متناظر در زیر دنباله‌های مشترک دقیقاً منطبق باشند، کار را کمی سخت می‌کند. این مشکل با استفاده از ترانسها ( $\Rightarrow$ ) هنگام مقایسه عناصر قابل حل است. بنابراین، مطابق گفته بالا، دو جزء  $a$  و  $b$  (به ترتیب از دنباله‌های  $X$  و  $Y$ ) مطابقتند اگر:

$$a(1-\varepsilon) < b < a(1+\varepsilon)$$

به جای ترانس نسبی، ممکن است از ترانس مطلق استفاده شود، یعنی  $b$  باید بین  $a - \varepsilon$  و  $a + \varepsilon$  باشد [۵]. برای آشنایی بیشتر با این الگوریتم، مثال زیر توضیح داده می‌شود.

مثال: فرض کنید دنباله‌های  $X$  و  $Y$  به صورت زیر داده شده باشند. برای به دست آوردن  $LCSS$  این دو دنباله، قدمهای زیر را دنبال می‌کنیم.

$$\begin{aligned} X &= ABCB \\ Y &= BDCAB \end{aligned}$$

	$Y_j$	B	D	C	A	B
$X_i$						
A						
B						
C						
B						

همه عناصر رشته‌ها در ستونها و ردیفها می‌چینیم

	$Y_j$	B	D	C	A	B
$X_i$	.	.	.	.	.	.
A	.					
B	.					
C	.					
B	.					

همه عناصر ستون و ردیف اول را صفر می‌کنیم

	$j$	۰	۱	۲	۳	۴	۵
	$Y_j$	(B)	D	C	A	B	
(A)	.	.	.	.	.	.	
B	.						
C	.						
B	.	.					

طبق قدمهای الگوریتم همه عناصر ستونها و ردیفها را با هم مقایسه کرده و عناصر متناظر در جدول را مقداردهی می‌کنیم اول را صفر می‌کنیم

	$Y_j$	B	D	C	A	B
$X_i$	.	.	.	.	.	.
A	.	.	.	.		
B	.					
C	.					
B	.					

	$Y_i$	B	D	C	A	B
$X_i$	.	.	.	.	.	.
A	.	.	.	.	1	
B	.					
C	.					
B	.					

	$Y_i$	B	D	C	A	B
$X_i$	.	.	.	.	.	.
A	.	.	.	.	1 → 1	
B	.					
C	.					
B	.					

	$Y_i$	B	D	C	A	B
$X_i$	.	.	.	.	.	.
A	.	.	.	.	1	1
B	.	→ 1	→ 1	→ 1	1	1
C	.					
B	.					

	$Y_i$	B	D	C	A	B
$X_i$	.	.	.	.	.	.
A	.	.	.	.	1	1
B	.	1	1	1	1	1
C	.					
B	.					

	$Y_i$	$B$	$D$	$C$	$A$	$B$
$X_i$	*	*	*	*	*	*
$A$	*	*	*	*	1	1
$B$	*	1	1	1	1	2
$C$	*	1	1	1		
$B$	*					

	$Y_i$	$B$	$D$	$C$	$A$	$B$
$X_i$	*	*	*	*	*	*
$A$	*	*	*	*	1	1
$B$	*	1	1	1	1	2
$C$	*	1	1	2		
$B$	*					

	$Y_i$	$B$	$D$	$C$	$A$	$B$
$X_i$	*	*	*	*	*	*
$A$	*	*	*	*	1	1
$B$	*	1	1	1	1	2
$C$	*	1	1	2	2	2
$B$	*					

	$Y_i$	$B$	$D$	$C$	$A$	$B$
$X_i$	*	*	*	*	*	*
$A$	*	*	*	*	1	1
$B$	*	1	1	1	1	2
$C$	*	1	1	2	2	2
$B$	*					

	$Y_j$	B	D	C	A	B
$X_i$	.	.	.	.	.	.
A	.	*	*	*	*	*
B	.	*	*	*	*	*
C	.	*	*	*	2	2
(B)	.	*	*	2	2	2

	$Y_j$	B	D	C	A	(B)
$X_i$	.	.	.	.	.	.
A	.	*	*	*	*	*
B	.	*	*	*	*	*
C	.	*	*	*	2	2
(B)	.	*	*	2	2	2

	$Y_j$	B	D	C	A	B
$X_i$	.	.	.	.	.	.
A	↑	↑	↑	↑	↖1	↖1
B	↑	↑	↑	↑	↖1	↖2
C	↑	↑	↑	↖1	↖1	↖2
B	↖1	↖1	↖1	↖2	↖2	↖3

$i$		$Y_j$	(B)	D	(C)	A	(B)
	$X_i$	.	.	.	.	.	.
1	A	*	*	*	*	*	*
2	(B)	*	↖1	↖1	*	*	*
3	(C)	*	*	*	↖2	↖2	*
4	(B)	*	*	*	↖2	↖2	↖3

LCS (reversed order): B C B

LCS (straight order): B C B

شکل ۱۶-۸) قدم‌های الگوریتم

### خوشه‌بندی جریان کلیکها<sup>۱</sup> با کمک بزرگترین دنباله‌های مشترک وزنی

گروه‌بندی بازدیدکنندگان براساس تعامل آنها با یک وب سایت، یک مشکل کلیدی در کاوش‌های کاربردی وب است. جریان کلیکهای تولید شده به‌وسیله کاربرهای مختلف اغلب الگوهای مشخصی را دنبال می‌کنند. برای خوشه‌بندی کاربران وب، براساس جریان کلیک در یک وب سایت و زمان صرف شده روی هر صفحه، از یک الگوریتم LCSS استفاده شده است [11].

با افزایش سریع کاربردهای تجارت الکترونیک، آگاهی از رفتار کاربر بر اساس تعاملاتش با یک وب سایت برای صاحبان وب سایت اهمیت بیشتری پیدا کرده است. تشخیص رفتار هر کاربر در بازدید از یک وب سایت می‌تواند مدیران سایت را برای تهیه محتوای سفارش شده برای دیگر کاربران قادر سازد. این موضوع کاربردهای زیادی در کسب وکار دارد. جریان کلیکهای یک کاربر، دنباله‌ای از صفحه‌های بازدید شده توسط او در یک وب سایت خاص در یک جلسه<sup>۲</sup> می‌باشد. هدف، خوشه‌بندی کاربران بر اساس جریان کلیکها در یک وب سایت خاص و یافتن گروههایی از کاربران است که با علایق و انگیزه‌های مشابه از یک سایت بازدید می‌کنند. در نتیجه همبستگی قوی بین جریان کلیکهای کاربران وجود دارد که نشان دهنده تشابه علایق کاربران است.

### ۸-۹- روشهای شاخص‌گذاری برای جستجوی تشابه در سریهای زمانی

مسئله دیگری که در بحثهای کاربردی سریهای زمانی مطرح است مسئله شاخص‌گذاری/بازیابی<sup>۳</sup> است [5]. مجموعه سریهای زمانی  $\{Y_1, \dots, Y_N\} = S$  را در نظر گرفته و سری مورد نظر  $X$  را داریم. سریهای زمانی موجود در  $S$  را که بیشترین تشابه با سری  $X$  دارند، پیدا می‌کنیم. برای مثال، روزهایی از سال که یک سهم مشخص تغییرات مشابهی مانند امروز داشته باشد را جستجو می‌کنیم.

<sup>۱</sup>- Clickstream

<sup>۲</sup>- Session

<sup>۳</sup>- Indexing/Retrieval

مسئله دیگر، شاخص‌گذاری زیر دنباله است. مجموعه دنباله‌های  $\mathcal{S}$  و دنباله یا الگوی مورد نظر  $X$  داده شده است. دنباله‌ای را در  $\mathcal{S}$  می‌یابیم که شامل زیر دنباله‌های مشابه  $X$  باشد. برای حل کارآی این نوع مسائل، باید از روشهای مناسب شاخص‌گذاری استفاده کنیم.

مسئله تشابه، مرتبط با مسئله شاخص‌گذاری می‌باشد. معمولاً تعیین شاخص برای روشهای اندازه‌گیری ساده شباهت، آسان و احتمالاً کم‌دققت است. ولی اندازه‌گیری شباهتهای پیچیده، تعیین شاخص را دشوار و جالب می‌کند.

یک سری زمانی با طول  $n$  می‌تواند به عنوان یک مشاهده در فضای  $n$  بعدی درنظر گرفته شود. شاخص‌گذاری مستقیم در این فضا به علت ابعاد خیلی زیاد آن ناکارآ می‌باشد. راه حل، استفاده از روشهای کاهش بعده است که در آن سری زمانی  $X$  با  $n$  مشاهده درنظر گرفته شده، چند مشخصه کلیدی آن استخراج و به نقطه  $f(X)$  در فضای مشخصه با ابعاد کمتر  $k$  نگاشت می‌شود (امید داریم که  $k < n$  باشد). این نگاشت باید طوری انجام شود که تشابه یا فاصله بین  $X$  و  $Y$  تقریباً با فاصله اقلیدسی دو نقطه  $f(X)$  و  $f(Y)$  برابر باشد. می‌توان از روشهای شناخته شده دسترسی فاصله‌ای برای شاخص‌گذاری فضای مشخصه‌های دارای ابعاد کمتر استفاده کرد، مانند

*R-trees*

.*VP-trees* یا *kd-trees*

در بسیاری از موارد سرعت دستیابی به سریهای زمانی مشابه با سری زمانی مورد جستجو، حائز اهمیت است. این مسئله در بسیاری از کاربردها دیده می‌شود. مثلاً یافتن سهامهایی که شبیه یک سهام خاص رفتار می‌کنند، پیدا کردن محصولاتی که چرخه تقاضای یکسانی دارند و پیدا کردن ژنهایی که الگوی میان آنها شبیه ژن خاصی است. این کاربردها نیازمند مکانیزم بازیابی برای دسته‌بندی سریهای زمانی با استفاده از روش دسته‌بندی نزدیک‌ترین همسایگی جهت بهبود زمان اجرا در الگوریتمهای خوبه‌بندی یا برای تحلیل اکتساب در داده‌های سریهای زمانی هستند.

شاخص‌گذاری در سریهای زمانی، توجه بسیاری از محققان را در سالهای اخیر جلب کرده است. با توجه به گسترش روزافزون اندازه بانکهای اطلاعاتی، شاخص‌گذاری می‌تواند در داده‌کاوی بسیار مؤثر باشد. اگر کاربر بخواهد در بانکهای اطلاعاتی گستره شروع به اکتشاف کند، باید داده‌ها به شکلی سازماندهی شوند که او بتواند به شکلی مؤثر و کارا داده‌های مورد

نظر خود را بازیابی کند. به طور عمومی می‌توان مسئله بازیابی سریهای زمانی را به صورت زیر تعریف کرد. بانک اطلاعاتی  $D$  شامل مجموعه‌ای از سریهای زمانی داده شده است. یک روش پیش‌پردازشی تعریف می‌کنیم که هدف آن پیدا کردن کارآی سری  $X$  عضو  $D$  نزدیک به سری داده شده  $Q$  است، (سری  $Q$  لزوماً در بانک اطلاعاتی وجود ندارد). برای حل این مسئله باید موارد زیر را در نظر گرفت:

- یکتابع فاصله که با درک کاربر از آنچه شباهت نامیده می‌شود، مطابقت دارد.
- یک رویه مؤثر شاخص‌گذاری که سرعت جستجوی کاربر را بالا می‌برد.

در زیربخش قبلی روش‌های مختلف برای تعریف شباهت (یا فاصله) بین دو سری زمانی بررسی شد. آسان‌ترین روش، تعریف فاصله بین دو سری با نگاشت هریک برروی یک بردار وسیس استفاده از نرم  $L_p$  برای محاسبه بود. فاصله نرم  $p$  بین دوبردار  $n$  بعدی  $\bar{x}$  و  $\bar{y}$  به صورت زیر تعریف می‌شود:

$$L_p(\bar{x}, \bar{y}) = \left( \sum_{i=1}^n (x_i - y_i)^p \right)^{\frac{1}{p}} \quad (16-8)$$

برای  $p=2$  این فرمول همان فاصله اقلیدسی مشهور و برای  $p=1$  فاصله مانهاتان است. اگر چه محاسبه چنین مقیاس فاصله‌ای آسان است، اما به کوچک‌ترین تغییرات در محور زمان بسیار حساس بوده و برای داده‌های مغشوش به خوبی عمل نمی‌کند. همان‌طور که گفته شد، انعطاف‌پذیرترین روشها برای تعریف تشابه در سریهای زمانی، روش‌های  $DTW$  و  $LCSS$  هستند.

صورت مسئله برای کاربران مختلف می‌تواند متفاوت باشد. مثلاً کاربران می‌توانند به دنبال تطابق کل دنباله و یا تنها به دنبال تطابق یک زیردنباله باشند. مقیاس‌های اندازه‌گیری تشابه یا فاصله، کاربردهای مختلفی دارند. همچنین کاربر ممکن است علاقه‌مند باشد که  $k$  تا از شیوه‌ترین سریهای زمانی که با فاصله  $\epsilon$  از سری مورد جستجو قرار دارند را بیابد. شاخص‌گذاری سریهای زمانی در دو حالت مورد بررسی قرار می‌گیرد که در حالت اول تابع فاصله متريک و در حالت بعدی تابع فاصله غير متريک است. در ادامه به طور خلاصه به بررسی هر یک از این حالات پرداخته می‌شود.

### شاخص گذاری سریهای زمانی با تابع فاصله متریک

شاخص گذاری علاوه بر اینکه موارد مشابه را درسازمان داده‌ها گردآوری می‌کند، امکان هرس داده‌های غیر مرتبط را نیز فراهم می‌نماید. هرس کردن کاملاً به متریک بودن تابع فاصله بستگی دارد.

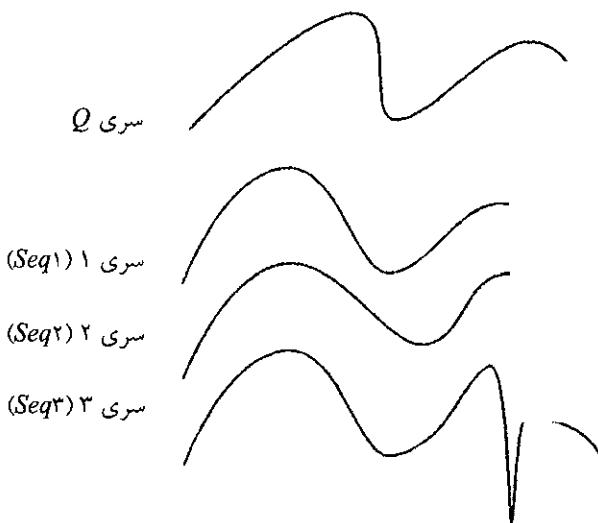
یک تابع فاصله  $d(X, Y)$  بین دو شیء  $X$  و  $Y$  متریک است اگر دارای شرایط زیر باشد:

$X=Y$  اگر  $d(X, Y) = 0$  و  $d(X, Y) \geq 0$  مثبت بودن

$d(X, Y) = d(Y, X)$  تقارن

$d(X, Y) + d(Y, Z) \geq d(X, Z)$  نامساوی مثلثی

فاصله‌اقلیدسی یک تابع فاصله متریک است. قدرت هرس کردن چنین تابعی در مثال زیر نشان داده شده است.



شکل ۱۷-۸) قدرت هرس ناساری مثلثی

مثال: فرض کنید مجموعه دنباله‌های  $S = \{Seq1, Seq2, Seq3\}$  در شکل ۱۷-۸) داده شده است. همچنین فرض کنید فاصله سه دنباله را با روش فواصل جفتی محاسبه و در جدول ۸-۱) مرتب کرده‌ایم.

جدول ۱-۸) فواصل جفنی

سری ۳	سری ۲	سری ۱	سری ۱
سری ۱	۲۰	۰	۱۱۰
سری ۲	۰	۲۰	۹۰
سری ۳	۹۰	۱۱۰	۰

برای پیدا کردن شبیه‌ترین دنباله به دنباله  $Q$  از بین سه دنباله فوق بهترین روش، تصویر دنباله‌هاست. در این روش فاصله همه آنها با  $Q$  محاسبه شده و یکی از آنها که کوتاه‌ترین فاصله را با  $Q$  دارد، انتخاب می‌شود. اگر تابع فاصله، نامساوی مثلثی را برآورده کند و داشته باشیم  $D(Q, Seq1) = 20$  و  $D(Q, Seq2) = 120$  آنگاه به این علت که:

$$D(Q, Seq3) \geq D(Q, Seq2) - D(Seq2, Seq3) \rightarrow D(Q, Seq3) \geq 120 - 90 = 40$$

ما می‌توانیم به راحتی  $Seq3$  را حذف کنیم زیرا راه حل بهتری را پیشنهاد نمی‌کند. برای بهبود بیشتر باید از یک شاخص چند بعدی استفاده کنیم.

شاخصهای چند بعدی و پیدا کردن نزدیک‌ترین همسایگی‌ها به شکلی مؤثر در فضایی با ابعاد بالا، توجه بسیاری از متخصصین را در علوم کامپیوتری و تحقیقات بانکهای اطلاعاتی به خود جلب کرده است. یک روش ساده شاخص‌گذاری، درنظر گرفتن سری زمانی با طول  $n$  به عنوان یک نقطه  $n$  بعدی است. ما می‌توانیم هر سری زمانی را به عنوان یک نقطه در ساختار  $n$  بعدی  $R\text{-tree}$  ذخیره کنیم. برای پیدا کردن نزدیک‌ترین همسایگی، با تبدیل هر سری زمانی به یک نقطه  $n$  بعدی و استفاده از ساختار شاخص‌گذاری مانند استفاده از  $R\text{-tree}$  به جستجوی نزدیک‌ترین همسایگی می‌پردازیم.

متاسفانه این ایده در عمل به خوبی کار نمی‌کند زیرا طول بلند سریهای زمانی معمولاً نقاطی با ابعاد بسیار بالا می‌سازد. هنگامی که ابعاد زیاد شود، کارآیی ساختارهای شاخص‌گذاری مختلف به تدریج کاهش می‌یابد.

**شاخص گذاری تشابه سریهای زمانی بازگشتی با تابع فاصله غیر متريک**  
توابع فاصله که نسبت به داده‌های بسیار مغلوش، مقاوم هستند، معمولاً نامساوی مثلثی را نقض می‌کنند. چنین توابعی، همه بخشها در سری زمانی را یکسان در نظر نمی‌گيرند. اگر چه

این رفتار مفید است، زیرا مدل صحیح‌تری از ادراک بشری را نشان می‌دهد. زمانی که مردم هر نوع داده‌ای (تصویر، سری زمانی و . . .) را مقایسه می‌کنند، بیشتر روی قسمتهایی که شبیه هستند، تمرکز کرده و توجه کمتری به قسمتهایی که شبیه نیستند، می‌کنند. فواصل غیرمتريک امروزه در بسیاری از دامنه‌ها به کار گرفته می‌شوند، مانند تطابق رشته (*DNA*), فیلتر کردن مشارکتی (هنگامی که مشتری با الگوی از پیش ذخیره شده مشتریان منطبق شود) و بازیابی تصاویر مشابه از بانکهای اطلاعاتی. علاوه بر آن تحقیقات علم روانشناسی مطرح می‌کند که قضاوت‌های مشابه بشری نیز غیرمتريک هستند. به علاوه برای توابع فاصله غیرمتريک، حالت‌های زیر مطرح می‌شود:

- سریهای زمانی در نرخ‌های نمونه‌گیری متفاوت یا سرعت‌های مختلف جمع‌آوری می‌شود. سریهای زمانی به دست آمده نتایج نمونه‌گیری در فواصل زمانی ثابت را تضمین نمی‌کنند. گیرنده‌های جمع‌آوری کننده داده، ممکن است برای یک دوره زمانی خاص، یکسان عمل کنند و به نرخ‌های نمونه گیری متناقض منتهی شوند. علاوه بر آن وقتی دوسری زمانی دقیقاً در یک جهت حرکت می‌کنند، ولی یکی با سرعتی دو برابر نسبت به دیگری در حرکت است، نتیجه به احتمال زیاد، فاصله‌اقلیدسی بسیار بزرگی است.
- سریهای زمانی شامل نقاط پرت هستند. نقاط پرت به دو صورت تعریف می‌شوند. یکی رخداد غیر عادی در برگیرنده جمع‌آوری کننده داده و دیگری واکنش به خطای بشری می‌باشد.
- سریهای زمانی طول‌های مختلفی دارند. فاصله‌اقلیدسی، مربوط به سریهای زمانی با طول یکسان است. هنگامی که طول‌ها متفاوت هستند، ما باید تصمیم بگیریم که آیا می‌خواهیم سری طولانی‌تر را کوتاه کنیم یا سری کوتاه‌تر را با جملات صفر لایه‌گذاری کنیم. مثالهایی از چنین مترهای فاصله، توابع فاصله *LCSS* و *DTW* می‌باشند. مشاهده اینکه هیچ‌کدام از آن دو متریک نیستند ساده است.

مثال: سه دنباله زیر را داریم:

$$Seq1 = (0, 1, 1, 0, 0, 1, 0, 0)$$

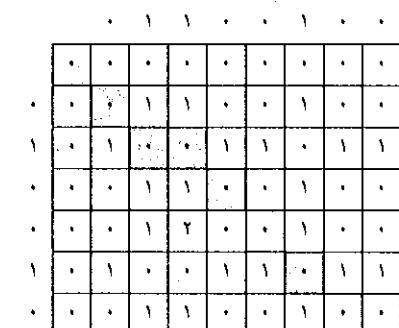
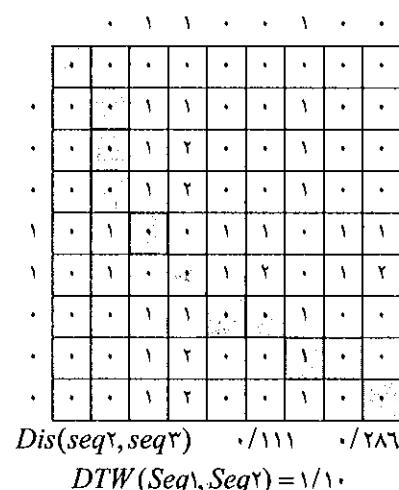
$$Seq2 = (0, 0, 0, 1, 1, 0, 0, 0)$$

$$Seq3 = (0, 0, 1, 0, 0, 1, 0)$$

$$DTW(Q, C) = \min \left\{ \sqrt{\sum_{k=1}^K w_k} / K \right\}$$

$$LCSS(X, Y) = (m + n - \text{Dis}) / (m + n)$$

	<i>DTW</i>	<i>LCSS</i>
(seq <sup>1</sup> , seq <sup>2</sup> ) <i>Dis</i>	۱/۰	۰/۲۵۰
(seq <sup>1</sup> , seq <sup>r</sup> ) <i>Dis</i>	.	۰/۱۴۳



.	.	.	.	۱	۱	۱	.	.
.	۱	۱	۱	۱	۱	۱	۱	۱
۱	۰	۱	۱	۱	۲	۲	۲	۲
۱	۰	۱	۱	۱	۲	۳	۳	۳
.	۰	۱	۲	۳	۳	۳	۴	۴
.	۰	۱	۲	۳	۳	۳	۴	۵
۱	۰	۱	۲	۳	۴	۴	۴	۰
.	۰	۱	۲	۳	۴	۴	۰	۰
۱	۰	۱	۲	۳	۴	۴	۰	۷
.	۰	۱	۲	۳	۴	۴	۰	۷

$$DTW(Seq^r, Seq^r) = 1/9$$

.	.	.	.	۱	۱	۱	.	.
.	۱	۱	۱	۱	۱	۱	۱	۱
۱	۰	۱	۱	۱	۲	۲	۲	۲
۱	۰	۱	۱	۱	۲	۳	۳	۳
.	۰	۱	۱	۱	۲	۳	۳	۳
.	۰	۱	۱	۱	۲	۳	۳	۴
۱	۰	۱	۱	۱	۲	۳	۳	۳
.	۰	۱	۱	۱	۲	۳	۳	۴
۱	۰	۱	۱	۱	۲	۳	۳	۴

$$LCSS(Seq^l, Seq^r) = (\lambda + \lambda - \gamma \times \gamma) / (\lambda + \lambda) = 1/\varepsilon$$

.	.	.	.	۱	۱	.	.	.
.	۱	۱	۱	۱	۱	۱	۱	۱
۱	۰	۱	۱	۲	۲	۲	۲	۲
۱	۰	۱	۱	۲	۳	۳	۳	۳
.	۰	۱	۱	۲	۳	۳	۴	۴
.	۰	۱	۱	۲	۳	۳	۴	۴
۱	۰	۱	۱	۲	۳	۳	۴	۰
.	۰	۱	۱	۲	۳	۳	۴	۰
۱	۰	۱	۱	۲	۳	۳	۴	۰

$$LCSS(Seq^r, Seq^l) = (\lambda + \gamma - \gamma \times \delta) / (\lambda + \gamma) = \varepsilon / 1\varepsilon$$

	۰	۱	۰	۰	۱	۰	۰	۰	۰	۰	۰
۰	۰	۱	۱	۱	۱	۱	۱	۱	۱	۱	۱
۱	۰	۱	۲	۲	۲	۲	۲	۲	۲	۲	۲
۱	۰	۱	۲	۲	۲	۲	۳	۳	۳	۳	۳
۰	۰	۱	۲	۳	۳	۳	۳	۴	۴	۴	۴
۰	۰	۱	۲	۳	۴	۴	۴	۴	۵	۵	۵
۱	۰	۱	۲	۳	۴	۵	۵	۵	۶	۶	۶
۰	۰	۱	۲	۳	۴	۵	۵	۶	۶	۶	۶
۰	۰	۱	۲	۳	۴	۵	۶	۶	۶	۶	۶

$$LCSS(Seq1, Seq3) = (8+6-2 \times 6) / (8+6) = 1/7$$

شکل ۸-۸) قدم‌های الگوریتم

$$DTW(Seq1, Seq2) = 1/9$$

$$LCSS(Seq1, Seq2) = (8+8-2 \times 6) / (8+8) = 1/4$$

$$LCSS(Seq3, Seq2) = (8+6-2 \times 5) / (8+6) = 4/14$$

$$LCSS(Seq1, Seq3) = (8+6-2 \times 6) / (8+6) = 1/7$$

برای شاخص گذاری مناسب سریهای زمانی براساس مفهوم تشابه، نیاز به استفاده از روش‌های کاهش بعد می‌باشد. در ادامه به بررسی ضرورت استفاده از روش‌های کاهش بعد و تبدیل داده‌ها پرداخته می‌شود.

### روشهای تبدیل و کاهش داده

به سبب اندازه بسیار زیاد و ابعاد بالای داده‌های به شکل سری‌های زمانی، روش‌های کاهش داده به عنوان اولین قدم در تحلیل سریهای زمانی به کار می‌رود. روش‌های کاهش داده نه تنها منجر به اشغال فضای کمتر می‌شود، بلکه سرعت پردازش داده‌ها را نیز افزایش می‌دهد. همان‌طور که در فصل دوم اشاره شد، مهم‌ترین استراتژی برای کاهش داده، انتخاب زیرمجموعه‌ای از ویژگیها می‌باشد که در آن ویژگیهای نامرتب و اضافی این داده‌ها در نظر گرفته نمی‌شود. علاوه بر این روش‌های کاهش بعد و کاهش حجم داده‌ها (نظریه نمونه گیری، خوشه‌بندی و...) نیز از دیگر ویژگیهای مؤثر و کارآمد برای کاهش مقادیر عظیم داده‌ها می‌باشد. بدلیل اینکه سریهای زمانی به عنوان یک نوع از داده‌های با حجم بالا مد نظر قرار می‌گیرند و هر نقطه بر حسب زمان

به عنوان یک بعد در نظر گرفته می‌شود، کاهش بعد یکی از مهم‌ترین مباحث مرتبط با تحلیل سریهای زمانی است. یکی از دلایل توجه روزافزون به کاهش بعد، در تحلیل و داده‌کاوی سریهای زمانی، کاهش پیچیدگی محاسبات در اثر کاهش حجم و بعد داده‌ها می‌باشد [۶].

ایده کلیدی شاخص گذاری مؤثر سریهای زمانی، کاهش بعد فضاست. برای کاهش ابعاد فضا، نمونه  $n$  بعدی را که نشان‌دهنده سری زمانی است، در یک فضای  $k$  بعدی ( $n < k$ ) تصویر می‌کنیم، به طوری که فاصله‌ها تاجای ممکن بدون تغییر باقی بمانند. سپس می‌توانیم از یک روش شاخص گذاری روی فضای جدید که ابعاد کمتری دارد، استفاده کنیم. چارچوب عمومی روش *GEMINI*، برای شاخص گذاری سریهای زمانی، با روش‌های کاهش بعد از گامهای زیر تبعیت می‌کند.

- مجموعه سریهای زمانی را بر روی یک فضای کاهش بعد یافته، نگاشت می‌کنیم.
- از یک روش شاخص گذاری برای شاخص گذاری فضای جدید، استفاده می‌کنیم.
- برای سری زمانی مورد جستجوی  $Q$ ، دنباله  $Q$  را بر روی فضای جدید نگاشت کرده و نزدیک‌ترین همسایگی‌ها به  $Q$  را در فضای جدید، با استفاده از شاخص گذاری پیدا می‌کنیم. سپس فواصل واقعی را محاسبه کرده و نزدیک‌ترین را انتخاب می‌کنیم.

روشهای کاهش بعد متعددی در تحلیل سریهای زمانی مورد استفاده قرار می‌گیرند که از آن جمله می‌توان به تبدیل فوریه گستته، تبدیل موجک گستته، تجزیه مقدار منفرد بر مبنای تحلیل مؤلفه‌های اصلی و تصویرکردن تصادفی اشاره کرد. این فنون در فصل دوم، قسمت کاهش بعد توضیح داده شده‌اند.

تبدیلات *SVD* دارای مزیت کاهش بعد بهینه تصاویر خطی می‌باشد. یعنی بهترین حفظ را از میانگین مریع خطا بین تصاویر اصلی و تصاویر تقریبی انجام می‌دهد. اگر چه محاسبه آن در مقایسه با روش‌های دیگر دشوار است، مخصوصاً اگر سریهای زمانی خیلی طولانی باشند. علاوه بر این، این روش برای شاخص گذاری زیر دنباله‌ها کاربرد ندارد.

تبدیلات گستته فوریه، طیف فرکانس یک سیگنال یک بعدی را توصیف می‌کند. روش *DFT* به عنوان یک روش کاهش بعد برای سریهای زمانی ارائه شده است.

تجزیه موجک<sup>۱</sup> که سریهای زمانی را به شکل مجموع توابع اولیه نشان می‌دهد، مشابه روش DFT می‌باشد. ساختار روش WD با DFT فرق دارد، زیرا تأثیر ضرایب مختلف در آن بیشتر در زمان متغیر کر شده‌اند تا در فرکانس. فواید WD آن است که سریهای زمانی تبدیل شده در همان دامنه (دامنه موقعت) باقی می‌ماند و الگوریتم کارآئی با پیچیدگی  $O(n)$  برای محاسبه تبدیلات وجود دارد. معایب آن، مشابه روش DFT است.

تصویر کردن تصادفی یک روش کاهش بعد عمومی است که در سال ۱۹۹۸ ارائه و در سال ۲۰۰۱ برای حوزه سریهای زمانی به کار گرفته شد.

استفاده از مقیاس‌بندی چند بعدی برای شاخص‌گذاری سریهای زمانی، دشوار است. روش MDS، نگاشت یک مجموعه از سریهای زمانی به نقاط  $k$  بعدی است (که  $k$  یک مقدار کوچک است). برای پاسخ‌گویی به جستجوی تشابه، باید قادر باشیم سری زمانی قابل جستجو را روی فضای  $k$  بعدی نگاشت کنیم. به علت نوع روشی که الگوریتم MDS برای یافتن جستجوهای مورد نظر استفاده می‌کند، باید فواصل همه سریهای زمانی داخل مجموعه آن را بیابیم و این عملیاتی خطی است.

نگاشت سریع، تخمینی بوده و روشی بسیار شبیه به روش مقیاس‌دهی چندبعدی است. برای بررسی مبسوط روش‌های کاهش بعد به فصل پیش‌پردازش، بخش کاهش بعد رجوع کنید.

## تخمین قطعه‌ای خط<sup>۲</sup>

به جای استفاده از DFT برای تقریب یک سری زمانی می‌توانیم از تقریب چندجمله‌ای استفاده کنیم. اگر چه با استفاده از یک روش خطی، یک روش عمومی برای تعریف مقیاس اندازه‌گیری غیراقلیدسی به دست می‌آید، ولی هیچ روش شاخص‌گذاری عمومی قابل قبولی شناخته نشده است. تحقیقات اخیر نشان می‌دهند که شاخص‌گذاری این سریهای زمانی در صورتی که سریهای زمانی با تابع ثابت قطعه‌ای تقریب زده شود، امکان‌پذیر است. ایده اصلی این روش، کاهش بعد از طریق تقسیم سریهای زمانی به  $k$  قطعه هماندازه می‌باشد. مقدار میانگین داده‌هایی که در هر قطعه می‌افتد محاسبه شده و برداری از این مقادیر، نمایش داده‌های کاهش

<sup>۱</sup>- WD

<sup>۲</sup>- Line Segment Approximation

بعد یافته می‌باشد. این نوع نمایش داده‌ها، تخمینهای خطی-قطعه‌وار نامیده می‌شود. این نوع ارائه، نسبت به روش  $DFT$  مزایایی دارد. مثلاً تبدیلات می‌تواند در زمانهای خطی اجرا شود. مهم‌تر از آن، این روش گسترهای از مقیاس‌های اندازه‌گیری، مانند نرم  $L_p$  گستته،  $DTW$  و جستجوی اقلیدسی وزندار را پشتیبانی می‌کنند. این روش مانند  $WD$  است، هنگامی که  $K$  ضریب ۲ داشته و میانگین برای تعریف زیردباهه استفاده شود.

## منابع

- ۱) غاطمی قمی م، ت. (۱۳۷۰) «پیش‌بینی و تجزیه و تحلیل سریهای زمانی»، نشردانش امروز.
  - ۲) نیرومند ح، بزرگنیا ا. (۱۳۷۲) "مقدمه‌ای بر تحلیل سریهای زمانی"، نشردانشگاه فردوسی مشهد.
  - ۳) نیرومند، حسینعلی، ۱۳۷۱، تجزیه و تحلیل سریهای زمانی، نشردانشگاه فردوسی مشهد.
  - ۴) ابریشمی، حمید، ۱۳۷۳، اقتصادستنجه کاربردی
- 5) Ye N. (2003) "*The hand book of data mining*"
- 6) Han. J, Kamber. M,(2006) "*data mining concepts and techniques*"
- 7) Chu K. W. ,Wong M. H. ,(1998) "*fast time series searching with scaling and shifting*"
- 8) Keogh. E, Pazzani M. (2000) "Scaling up dynamic time warping for datamining application",*proceeding of the sixth acm sigkdd conference on knowledge discovery and data mining*
- 9) Keogh E. , Pazzani. M,(2001) "derivative dynamic time warping"
- 10) Keogh E. et al (2004) "Exact indexing of dynamic time warping"
- 11) Banerjee A. , Ghosh J. (2000) "clickstream clustering using weighted longest common subsequences"
- 12) Faloutsos C. , Lin K. (1995) "Fast map"
- 13) Smith L. I. (2002) "A tutorial on principal components analysis"

---

## فصل نهم

---

# تحلیل شبکه‌های اجتماعی

در دهه‌های اخیر، نظریه شبکه‌های اجتماعی (که در آن رابطه میان موجودیت‌ها<sup>۱</sup> به صورت پیوندهای درون یک گراف نشان داده می‌شوند) توجهات روز افزونی را به خود جلب کرده است. بدین ترتیب تحلیل شبکه‌های اجتماعی از دید داده‌کاوی، تحلیل پیوندها یا پیوندکاوی نیز نامیده می‌شود. در این بخش، ما ابتدا مفهوم شبکه‌های اجتماعی را مطرح نموده و به مطالعه ویژگی‌های شبکه‌های اجتماعی می‌پردازیم. سپس نگاهی به وظایف و چالشهای موجود در پیوندکاوی خواهیم داشت و در نهایت چند نمونه از شبکه‌های اجتماعی را مورد کاوش قرار می‌دهیم.

## ۱-۹- تعریف شبکه اجتماعی

از دید داده‌کاوی، شبکه اجتماعی مجموعه داده‌های ناهمگن<sup>۲</sup> (نامتجانس) و چند رابطه‌ای<sup>۳</sup> است که توسط یک گراف نمایش داده می‌شود. نوعاً گرافها بسیار بزرگ و مشکل از گره‌هایی معادل اشیاء و یالهایی<sup>۴</sup> معادل پیوندها (که معرف رابطه یا تعامل بین اشیا می‌باشند) هستند. هم

---

<sup>۱</sup>- Entity

<sup>۲</sup>- Heterogeneous

<sup>۳</sup>- Multirelational

<sup>۴</sup>- Edges

گره‌ها و هم پیوندها ویژگیهای دارند. اشیا ممکن است بر جسب دسته داشته باشند. پیوندها می‌توانند یک طرفه باشند و نزومی ندارد که دو حالته باشند.

لازم نیست شبکه‌های اجتماعی زمینه‌ای اجتماعی داشته باشند. مثالهای واقعی بسیاری از شبکه‌های اجتماعی تکنولوژیکی، تجاری، اقتصادی و زیستی وجود دارد؛ شبکه‌های توزیع نیروی الکتریسیته، گرافهای تماسهای تلفنی، انتشار ویروسهای کامپیوتری، شبکه گسترده جهانی، شبکه‌های استناد و هم‌نویسندهای دانشمندان. مثال دیگر عبارت است از شبکه‌های مشتریان در موقعی که توصیه محصول بر اساس اولویت سایر مشتریان صورت می‌گیرد. در بیولوژی، نمونه‌ها طیف وسیعی را از شبکه‌های شیوع بیماری، شبکه‌های متابولیک و سلولی، شبکه غذایی تا شبکه عصبی کرم (تنها موجودی که شبکه عصبی اش کاملاً نگاشت شده است) دربرمی‌گیرند. تبادل پیامهای پست الکترونیک بین شرکتها، گروه‌های خبری، اتفاقهای گپ، شبکه دوستی، هم‌پوشانی هیئت مدیره‌های شرکتها بزرگ آمریکایی و غیره نمونه‌هایی از حوزه جامعه‌شناسی هستند.

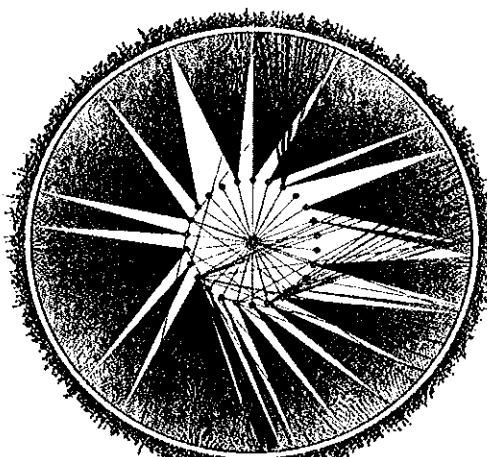
«شبکه‌های (اجتماعی) دنیای کوچک»<sup>۱</sup> اخیراً نظرات بسیاری را به خود معطوف داشته است. آنها مفهوم «دنیای کوچک» را منعکس می‌نمایند که در اصل بر شبکه‌های بین افراد مرکز دارد. همان عبارتی که نشان دهنده تعجب زیاد اولیه میان دو غریبه است، وقتی در می‌یابند که به صورت غیر مستقیم از طریق نفر سومی که هر دو با وی سابقه آشنا بی دارند، با هم رابطه دارند: «چه دنیای کوچکی!». در ۱۹۶۷ استنلی میلگرام<sup>۲</sup> جامعه‌شناس هاروارد و همکارانش آزمایشی را ترتیب دادند که در آن از مردم کانزاس و نبراسکا درخواست شده بود تا نامه‌هایی را به دست افرادی غریبه در بوستون برسانند، بدین صورت که این نامه‌ها را به دوستانی ارسال کنند که گمان می‌کنند ممکن است آن افراد غریبه را بشناسند. نیمی از نامه‌ها به‌طور موفقیت‌آمیزی از طریق کمتر از ۵ واسطه به مقصد رسیدند. مطالعات دیگری که میلگرام و سایرین در شهرهای دیگر صورت دادند، نشان داد که ظاهراً به‌طور عمومی «شش سطح جدایی»<sup>۳</sup> بین هر دو نفر در جهان وجود دارد. نمونه‌هایی از شبکه‌های دنیای کوچک در شکل (۱-۹) نشان داده شده است.

<sup>۱</sup>- Small World

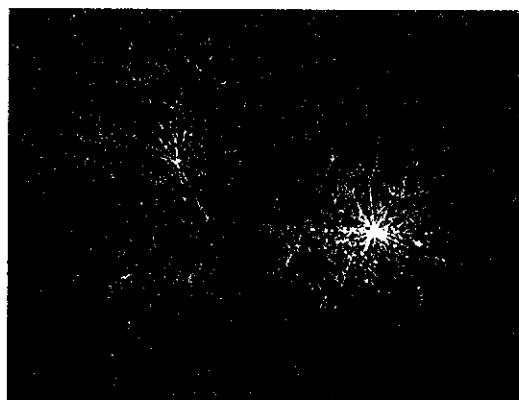
<sup>۲</sup>- Stanley Milgram

<sup>۳</sup>- Six Degrees of Separation

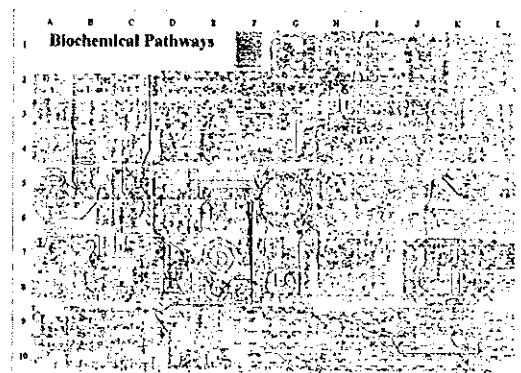
شبکه‌های دنیای کوچک با ویژگی دارا بودن درجه بالایی از خوشبذری محلی برای نسبت کوچکی از گره‌ها شناخته می‌شوند (به عنوان مثال این گره‌ها با دیگری به هم وصل می‌شوند)، که در عین حال بیش از چند سطح از سایر گره‌ها باقیمانده جدا نیستند. این باور وجود دارد که بسیاری از شبکه‌های زیستی، اجتماعی، و فیزیکی ساخته دست بشر این ویژگی‌های دنیای کوچک را به معرض نمایش می‌گذارند. این ویژگیها بعداً توضیح داده شده و مدلسازی می‌شوند.



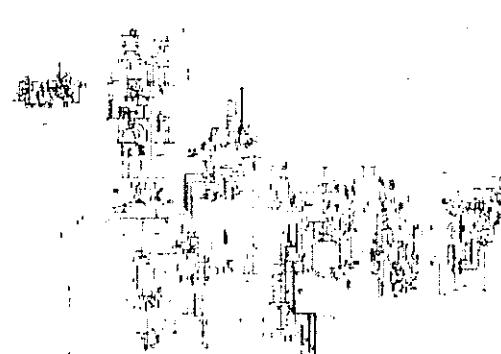
(a)



(b)



(c)



(d)

شکل ۱-۹) مثالهای شبکه اجتماعی در دنیای حقیقی: (a) هم‌نویستگی علمی، (b) صفحات مرتبط در بخشی از اینترنت، (c) مسیر بیوشیمی، (d) شبکه قدرت الکتریسیته نیویورک

«به طور کلی چرا این همه علاقه به شبکه‌های دنیای کوچک و شبکه‌های اجتماعی وجود دارد؟ چه فایده‌ای در مشخص نمودن ویژگی‌های شبکه‌ها و کاوش آنها به منظور بیشتر آموختن از ساختارشان وجود دارد؟» علت آن است که ساختار همواره بر عملکرد مؤثر است. به عنوان مثال، توبولوژی شبکه‌های اجتماعی بر شیوه بیماریهای مسری در یک جمعیت ساختار یافته تأثیر دارد. توبولوژی شبکه نیرو بر ثبات و انسجام انتقال نیرو مؤثر است.

به طور نمونه، یک نارسایی مرتبط در تاریخ ۱۴ آگوست ۲۰۰۳ در کلیولند واقع در اوهایو در سیستم شبکه ایجاد شده. این نارسایی منجر به از کار افتادن کارخانجات نیروی هسته‌ای در ایالت نیویورک و اوهایو گردید و باعث خاموشی گسترده در بخش‌های زیادی از شمال شرقی ایالات متحده و جنوب شرقی کانادا شد که حدود ۵۰ میلیون نفر را در بر می‌گرفت.

توجه به شبکه‌ها بخشی از مطالعات وسیع‌تری است که در توصیف کامل و دقیق «سیستمهای پیچیده<sup>۱</sup>» انجام می‌شود. قبل از شبکه‌هایی که برای مطالعات تجربی در دسترس بودند، محدود و کوچک بودند و اطلاعات ناچیزی از گره‌های آنها وجود داشت. باید از اینترنت ممنون باشیم که در حال حاضر مقادیر عظیمی از داده‌های مرتبط با شبکه‌های اجتماعی بسیار بزرگ را در دسترس ما قرار داده است. این شبکه‌ها نوعاً دهها هزار تا میلیون‌ها گره را دربرمی‌گیرد و اغلب اطلاعات وسیعی در سطح هر گره موجود می‌باشد.

دسترسی به کامپیوترهای قدرتمند نیز بررسی ساختار شبکه‌ها را ممکن ساخته است. مطالعه شبکه‌های اجتماعی می‌تواند به ما در دسترسی بهتر و آسانتر به سایر مردم دنیا کمک کند. بعلاوه، مطالعه دنیای کوچک، با توجه به جدایی نسبتاً کم بین گره‌هایشان، می‌تواند به ما در طراحی شبکه‌هایی که انتقال کارای اطلاعات یا سایر منابع را تسهیل می‌کنند، یاری رساند، بدون آنکه مجبور باشیم از شبکه‌ای با تعداد زیادی رابطه زائد استفاده کنیم. به عنوان مثال می‌تواند به ما در طراحی موتورهای جستجو هوشمندتری در وب کمک کند، به طوریکه در پاسخ به یک پرس‌وجو، وب سایتها مرتبطی که کمترین میزان جدایی از وب سایت اولیه را دارند، بیابند.

## ۲-۹- ویژگیهای شبکه‌های اجتماعی

همان‌طورکه در قسمت قبل توضیح داده شد، دانستن ویژگیهای شبکه‌های دنیای کوچک در موقعیتهای بسیاری مفید است. می‌توان مدل‌های مولد گراف را که دارای همین ویژگیها باشند، ایجاد نمود. تا قادر به پاسخگویی به سوالات «اگر – آنگاه» و پیش‌بینی چگونگی یک شبکه در آینده باشند. به عنوان مثال در مورد اینترنت، می‌پرسیم «اگر تعداد گره‌ها در اینترنت دو برابر شود، آنگاه اینترنت چگونه به نظر می‌رسد؟» و «تعداد يالها چندتا خواهد بود؟». چنانچه فرضیه‌ای با ویژگیهای عموماً پذیرفته شده، تناقض داشته باشد، پرچمی در برابر مشکوک بودن فرضیه بر می‌افرازد.

این موارد به کشف ناهنجاریها در گرافهای موجود کمک می‌نماید که ممکن است تقلب، هرزنگاری<sup>۱</sup> یا حمله‌های Ddos را مشخص کند. به علاوه مدل‌های مولد گراف می‌تواند برای شیوه سازی مواردی که گرافهای واقعی بیش از اندازه بزرگ هستند و بدین لحاظ جمع‌آوری‌شان غیرممکن است (مثل شبکه بسیار بزرگی از رابطه دوستی) کمک کند. در این قسمت، ویژگیهای اساسی شبکه‌های اجتماعی به همراه مدلی برای تولید گراف بررسی می‌شود. چه ویژگیهای شبکه‌های اجتماعی را مشخص می‌کند؟ اکثر مطالعات، امتیاز گره‌ها<sup>۲</sup> را بررسی کرده‌اند که همان تعداد يالهای متنه‌ی به هر گره است و یا فاصله بین هر دو گره را که با محاسبه طول کوتاهترین مسیر تعیین می‌شود. سایر محاسبات فاصله میان دو گره شامل فاصله متوسط<sup>۳</sup> و قطر مؤثر<sup>۴</sup> (به عنوان مثال، حداقل فاصله  $d$ ) به طوری که لاقل برای  $90\%$  گره‌ها در دسترس، طول مسیر حداقل  $d$  باشد) می‌باشد.

شبکه‌های اجتماعی به ندرت ایستا هستند. نمایش گراف آنها با افزایش یا حذف گره‌ها و يالها در طول زمان، تکامل می‌یابد. به‌طورکلی، شبکه‌های اجتماعی پدیده‌های زیر را نشان می‌دهند:

<sup>۱</sup>- Spam

<sup>۲</sup>- Nodes Degrees

<sup>۳</sup>- Average Distance

<sup>۴</sup>- Effective Diameter

قانون تراکم توانی<sup>۱</sup>: در گذشته، باور بر این بود که با تکامل یک شبکه، امتیاز گره‌ها به صورت خطی نسبت به تعداد گره‌ها افزایش می‌باید. این باور با نام فرض میانگین امتیاز ثابت<sup>۲</sup> شناخته می‌شد. با این وجود، آزمایشات گستردۀای نشان داده‌اند که بر عکس، با افزایش متوسط امتیاز در طول زمان، شبکه‌ها به طور فزاینده‌ای چگال می‌شوند (بدین ترتیب، تعداد یال‌ها به نسبت تعداد گره‌ها به صورت فوق خطی<sup>۳</sup> افزایش می‌یابند). این متراکم شدن از قانون تراکم توانی (یا قانون رشد توانی<sup>۴</sup>) پیروی می‌کند و عبارتست از:

$$e(t) \propto n(t)^a \quad (1-9)$$

که در آن  $e(t)$  و  $n(t)$  به ترتیب نمایانگر تعداد یال‌ها و گره‌های گراف در زمان  $t$  هستند و توان  $a$  عموماً بین ۱ و ۲ قرار می‌گیرد. توجه داشته باشید که چنانچه  $a = 1$  باشد، معادل میانگین امتیاز ثابت در طول زمان است، در حالی که  $a = 2$  بیانگر گرافی به شدت متراکم است که در آن هر گره توسط یال‌هایش به نسبت ثابتی از کل گره‌ها، متصل است.

جمع شدن قطر<sup>۵</sup>: به طور تجربی اثبات شده که قطر مؤثر با رشد شبکه کاهش می‌یابد. این با باور اولیه‌ای که معتقد بود قطر به آهستگی به صورت تابعی از اندازه شبکه افزایش می‌یابد، در تناقض است.

برای درک شهودی این مطلب یک شبکه استناد را در نظر بگیرید که در آن گره‌ها مقالات هستند و استناد (ارجاع) از یک مقاله به مقاله‌ای دیگر با یک یال جهت‌دار نمایش داده می‌شود. یال‌های خروجی از گره ۷ (نمایانگر مقالاتی که ۷ به آنها ارجاع داده است) در لحظه پیوستن به گراف «منجمد»<sup>۶</sup> هستند. متعاقباً به نظر می‌رسد فاصله در حال کاهش بین دو گره، مقالات بعدی باشد که به عنوان «پل»<sup>۷</sup> عمل کرده و مقالات اولیه را مورد ارجاع قرار داده‌اند.

<sup>۱</sup>- Densification Power Law

<sup>۲</sup>- Constant Average Degree Assumption

<sup>۳</sup>- Superlinearly

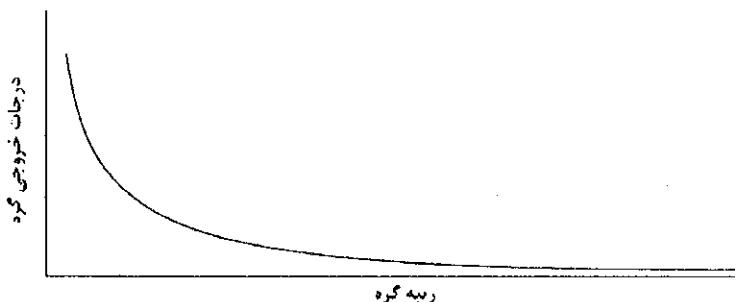
<sup>۴</sup>- Growth Power Law

<sup>۵</sup>- Shrinking Diameter

<sup>۶</sup>- Frozen

<sup>۷</sup>- Bridge

توزیع‌های دمپهن<sup>۱</sup> درجات ورودی و خروجی: با توجه به «قانون توان» تعداد امتیازات خروجی یک گره به پیروی از توزیعی دمپهن گرایش دارد:  $\frac{1}{n^a}$  که در آن  $n$  رتبه گره در ترتیب کاهش امتیاز خروجی است و نوعاً  $a < 2$  است (شکل ۲-۹). هرچه مقدار  $a$  کوچک‌تر باشد، دنباله پهن‌تر خواهد بود. این حالت بیانگر الحال ترجیحی<sup>۲</sup> است که در آن هر گره جدید با تعداد ثابتی یالهای خروجی به شبکه موجود وصل می‌شود و از قاعده «ثروتمند، ثروتمندتر می‌شود»<sup>۳</sup> پیروی می‌کند. امتیازات ورودی هم از توزیع دمپهن پیروی می‌کند که در ضمن نسبت به توزیع امتیازات خروجی چولگی بیشتری دارد.



شکل ۹ - ۲) تعداد درجات خروجی (محور y) یک گره تعابیل به پیروی از توزیعی دمپهن دارد. رتبه گره (محور x) ترتیب نزولی درجات خروجی گره است.

برای تولید گراف مدل آتش جنگل<sup>۴</sup> پیشنهاد شد که این خصوصیات تکامل گراف در طول زمان را دارا است. این مدل بر پایه این نظریه استوار است که گره‌های جدیدی از طریق سوزاندن<sup>۵</sup> یالهای موجود به طریق مسری به شبکه متصل شده و از دو پارامتر بهره می‌گیرد: احتمال سوزاندن پیشرو (p) و نسبت سوزاندن پسرو (r) که در ادامه توضیح داده می‌شوند. فرض کنید گره جدید  $t$  در زمان  $t$  اضافه شود، این گره در طی گامهای زیر به  $G_t$  گرافی که تاکنون ایجاد شده، متصل می‌شود:

<sup>۱</sup>- Heavy-Tailed

<sup>۲</sup>- Preferential Attachment Model

<sup>۳</sup>- Rich-Get-Richer

<sup>۴</sup>- Forest fire Model

<sup>۵</sup>- Burning

**گام اول:** به طور تصادفی یک گره سفیر<sup>۱</sup>،  $w$  انتخاب می‌کند و یک یال به  $w$  متصل می‌کند.

**گام دوم:**  $x$  یال متصل به  $w$  را انتخاب می‌کند.  $x$  مقداری تصادفی است که دارای توزیع برنولی با میانگین  $(p-1)$  است. هم یالهای ورودی به  $w$  و هم یالهای خروجی از  $w$  در نظر گرفته می‌شوند ولی یالهای ورودی را با احتمال  $r$  بار کمتر از یالهای خروجی برمی‌گزیند. گره‌هایی را که در سر دیگر یالهای انتخاب شده قرار دارند،  $w_x, w_1, \dots, w_n$  بنامید.

**گام سوم:** گره جدید ( $w$ ) یالهای خروجی را به  $w_x, w_1, \dots, w_n$  متصل نمود، حال گام ۲ را به طور بازگشتی برای هر کدام از گره‌های  $w_x, w_1, \dots, w_n$  انجام دهید. به منظور جلوگیری از افتادن در حلقه، هیچ گره‌ای برای بار دوم نباید انتخاب شود. این فرآیند تا زمانی که خاموش شود ادامه می‌یابد.

برای درک شهودی از مدل، به مثال خودمان از شبکه استناد، باز می‌گردیم. نویسنده مقاله جدید،  $v$ ، ابتدا به  $w$  مراجعه می‌کند. سپس یک زیر مجموعه از منابع  $w$  را دنبال می‌کند (که ممکن است پیشرو یا پسرو باشند) و به مقالات  $w_x, w_1, \dots, w_n$  دست می‌یابد. با ارجاع دادن این مقالات، تجمعیع مراجع به صورت بازگشتی ادامه می‌یابد.

بسیاری از مدل‌های تکامل شبکه بر گرافهای ایستا مبتنی است که ویژگیهای شبکه را بر اساس یک یا تعداد محدودی تصویر آنی از آن با کمی تأکید بر یافتن روندها در طول زمان تعیین می‌کنند. مدل آتش جنگل درحالی که به تکامل شبکه در طول زمان توجه دارد، ماهیت بسیاری از مدل‌های قبلی را درهم می‌آمیزد. مثلاً خاصیت امتیازهای خروجی دمپهن که به طبیعت بازگشتی تشکیل یالها تعلق دارد را نیز در نظر می‌گیرد؛ بدین ترتیب، گره‌های جدید شанс خوبی در سوزاندن بسیاری یالها و لذا ایجاد امتیازات خروجی بزرگ دارند. خاصیت امتیازات ورودی دمپهن نیز حفظ می‌شود، زیرا آتش جنگل از قاعده «ثروتمند، ثروتمندتر می‌شود» پیروی می‌کند. گره‌هایی که خیلی زیاد به هم متصل هستند، بدون در نظر گرفتن اینکه گره جدید از کدام گره سفیر آغاز می‌شود، به آسانی به گره‌ای جدید متصل می‌شوند. گونه‌ای از مدل کپی‌کننده<sup>۲</sup> نیز مورد نظر قرار گرفته است بدین صورت که یک گره جدید بسیاری از همسایگان گره سفیر خود

<sup>۱</sup>- Ambassador

<sup>۲</sup>- Copying Model

را کپی می‌کند. قانون تراکم توانی نیز تأیید می‌شود: یک گره جدید بالهای بسیاری در نزدیکی اجتماع گره سفیر خود خواهد داشت. مطالعات تجربی خاصیت جمع شدن قطر را تأیید می‌کنند. گره‌هایی که امتیازات خروجی دمپهن دارند، ممکن است به عنوان «پلهایی» به کار روند که قسمتهایی از شبکه را که پیش‌تر ناهمگن بودند، متصل کنند و قطر شبکه را کاهش دهند.

### ۹-۳- پیوندکاوی<sup>۱</sup>: وظایف و چالشها

چگونه می‌توان شبکه‌های اجتماعی را مورد کاوش قرار داد؟ روش‌های سنتی یادگیری ماشینی و داده‌کاوی که به عنوان ورودی، نمونه‌های تصادفی اشیاء همگنی را از یک رابطه واحد اخذ می‌کنند، ممکن است برای این حالت مناسب نباشد. داده‌های تشکیل‌دهنده شبکه‌های اجتماعی بیشتر ناهمگن، چندرابطه‌ای و نیمه‌ساخت‌یافته هستند. در نتیجه، حوزه جدیدی از پژوهش با نام پیوندکاوی پدید آمده است. پیوندکاوی محل تلاقی پژوهش در شبکه‌های اجتماعی، تحلیل پیوندها، وب‌کاوی و ابرمنتها<sup>۲</sup>، گراف‌کاوی، یادگیری رابطه‌ای<sup>۳</sup> و برنامه‌ریزی منطقی قیاسی<sup>۴</sup> است. پیوندکاوی متنضم مدل‌های توصیفی و پیش‌بینانه است. با درنظر گرفتن پیوندها (رابطه میان اشیاء) اطلاعات بیشتری برای فرآیند کاوش حاصل می‌شود. این امر وظایف جدید بسیاری را به همراه می‌آورد، لیست این وظایف با مثال‌هایی از حوزه‌های مختلف در زیر فهرست شده‌اند.

**دسته‌بندی اشیاء مبتنی بر پیوندها:** در روش‌های سنتی دسته‌بندی، اشیاء بر اساس ویژگی‌هایی که آنها را توصیف می‌کردند، دسته‌بندی می‌شدند. در دسته‌بندی مبتنی بر پیوندها، دسته یک شیء نه تنها بر اساس ویژگی‌ایش، بلکه بر اساس پیوند‌هایش و ویژگی‌های اشیایی که با آنها پیوند دارد، پیش‌بینی می‌شود.

**صفحات وب**، دسته یک صفحه هم بر اساس رخداد واژه‌ها (کلماتی که بر روی آن صفحه واقع

<sup>۱</sup>- Link Mining

<sup>۲</sup>- Hypertext

<sup>۳</sup>- Relational Learning

<sup>۴</sup>- Inductive Logic Programming

شده‌اند) و هم بر اساس متن لنگر<sup>۱</sup> (کلمات ابرپیوندی<sup>۲</sup> یعنی کلماتی که شما هنگام کلیک بر روی یک پیوند، بر روی آنها کلیک می‌کنید) انجام می‌شود و هر دوی آنها به عنوان ویژگیهای صفحه به کار می‌روند. به علاوه، دسته‌بندی بر پایه پیوندهای بین صفحات و سایر ویژگیهای صفحات و پیوندها نیز می‌باشد. در حوزه کتاب‌شناسی<sup>۳</sup>، اشیاء ما مقالات، نویسنده‌گان، مؤسسات، مجلات و کنفرانسها هستند. در این حالت وظیفه دسته‌بندی، پیش‌بینی مبحث یک مقاله است که بر اساس رخداد واژه‌ها، مورد ارجاع بودن (سایر مقالاتی که به این مقاله ارجاع داده‌اند)، و استناد کردن (سایر مقالاتی که در این مقاله به آنها استناد شده است) انجام می‌شود. استنادها به عنوان پیوند عمل می‌کنند. در مبحث شناخت بیماریهای همه‌گیر، وظیفه دسته‌بندی، پیش‌بینی نوع بیماری یک فرد است که بر اساس ویژگیهای (علائم بیماری) فرد بیمار و ویژگیهای سایر افرادی که فرد بیمار با آنها تماس داشته است (از این افراد با عنوان تماسهای بیمار یاد می‌شود) صورت می‌گیرد.

پیش‌بینی نوع شیء: این پیش‌بینی، نوع یک شیء را بر اساس ویژگیهای خودش و پیوندهایش و ویژگیهای اشیابی که به آن متصلند تعیین می‌کند. در حوزه کتاب‌شناسی ممکن است بخواهیم محل یک مطلب منتشر شده را مشخص کنیم مثلًاً کنفرانس یا مجله یا کارگاه. در حوزه ارتباطات، وظیفه مشابهی وجود دارد که پیش‌بینی آن است که آیا یک رابطه از طریق پست الکترونیک است، یا تلفن یا پست؟

پیش‌بینی نوع پیوند: این پیش‌بینی، بر اساس مشخصه‌های اشیاء درگیر در یک پیوند، نوع یا هدف آن پیوند را معلوم می‌سازد. به عنوان مثال با در دست داشتن داده‌ها شیوع بیماریهای همه‌گیر، ممکن است بخواهیم پیش‌بینی کنیم آیا دو نفر که یکدیگر را می‌شناسند از اعضای یک خانواده‌اند، یا همکارند، یا صرفاً با یکدیگر آشنا هستند. در نمونه‌ای دیگر ممکن است بخواهیم بدانیم آیا رابطه بین دو مؤلف به صورت مشاور-مشورت کننده است؟ با دستیابی به داده‌های صفحات وب می‌توانیم پیش‌بینی کنیم آیا یک پیوند روی یک صفحه، پیوند تبلیغاتی است یا پیوندی مربوط به پیمایش وب؟

<sup>۱</sup>- Anchor Text

<sup>۲</sup>- Hyperlink

<sup>۳</sup>- Bibliography

پیش‌بینی وجود پیوند: در پیش‌بینی نوع پیوند ما می‌دانستیم پیوندی بین دو شیء موجود است و می‌خواستیم نوع این پیوند را معلوم کنیم اما در این حالت می‌خواهیم پیش‌بینی کنیم آیا اصلًاً بین دو شیء پیوندی وجود دارد؟ مثال: پیش‌بینی اینکه آیا پیوندی بین دو صفحه وب وجود خواهد داشت؟ آیا مقاله‌ای به مقاله دیگری استناد خواهد کرد؟ و یا در علم شناخت بیماری‌های همه‌گیر می‌توانیم پیش‌بینی کنیم یک بیمار با چه کسی در تماس بوده است.

تخمین عدد اصلی<sup>۱</sup> پیوند: دو شکل تخمین عدد اصلی پیوند وجود دارد. اول، می‌توان تعداد پیوندهای متصل به یک شیء را پیش‌بینی نمود. به طور مثال تعیین مأخذیت<sup>۲</sup> یک صفحه وب می‌تواند بر اساس تعداد پیوندهایی که به آن متنه شده‌اند (پیوندهای ورودی) مشخص شود. به طور مشابه تعیین تعداد پیوندهای خروجی<sup>۳</sup> می‌تواند میان این باشد که آیا آن صفحه وب به عنوان قطب<sup>۴</sup> عمل می‌کند (منظور از قطب، یک یا مجموعه‌ای از صفحات وب است که به تعدادی صفحات مأخذ در همان مبحث استناد می‌کنند)، در حوزه کتاب‌شناسی، تعداد استنادات به یک مقاله می‌تواند نشانگر تأثیر آن مقاله باشد (هرچه استنادات به مقاله بیشتر باشد، انتظار می‌رود تأثیرگذاری بیشتری داشته باشد)، در علم بیماری‌های همه‌گیر، پیش‌بینی تعداد پیوندهای بین یک بیمار و تماس‌هایش، نشانگر انتقالهای بالقوه بیماری است.

حالت مشکل‌تر تخمین عدد اصلی پیوند، پیش‌بینی تعداد اشیایی است که در خلال یک مسیر از یک شیء قرار دارند. این مسئله در تخمین تعداد اشیائی که در پاسخ به یک پرس‌وجو باز می‌گردد، حائز اهمیت است. در حوزه صفحات وب، می‌توان تعداد صفحاتی را که از طریق پویش یک سایت بازیابی می‌شود، پیش‌بینی نمود (منظور از پویش، جستجویی خودکار و روشنمند در وب است که اساساً به منظور تهیه یک کپی از تمام صفحات مشاهده شده است تا برای پردازش‌های بعدی یک موتور جستجو استفاده شود). با در نظر گرفتن مسئله استناددهی، می‌توان از تخمین کاردینالیتی پیوندی برای پیش‌بینی تعداد استنادات یک نویسنده خاص در یک مجله خاص بهره برد.

<sup>۱</sup>- Cardinality

<sup>۲</sup>- Authoritativeness

<sup>۳</sup>- Out-Links

<sup>۴</sup>- Hub

**مصالحه<sup>۱</sup> اشیاء:** در مصالحه اشیاء، وظیفه پیش‌بینی این است که آیا در واقع دو شیء بر اساس ویژگیها و پیوندهایشان یکسان هستند؟ این وظیفه در استخراج اطلاعات<sup>۲</sup>، حذف دوباره‌کاری<sup>۳</sup>، یکی‌سازی اشیاء و انطباق<sup>۴</sup> استنادات امری معمول است و با عنوان عدم قطعیت موجودیت‌ها<sup>۵</sup> یا ارتباط رکورد<sup>۶</sup> نیز شناخته می‌شود. به عنوان نمونه: پیش‌بینی اینکه آیا دو وب سایت آینه یکدیگرند، آیا دو استناد واقعاً به یک مقاله استناد می‌کنند، و آیا دو درد آشکار حاصل از بیماری، در واقع یکی هستند؟

**کشف گروه‌ها:** کشف گروه نوعی خوشبندی است و به پیش‌بینی این امر می‌پردازد که چه هنگام مجموعه‌ای از اشیاء به یک گروه یا خوشه تعلق دارند و این کار را بر اساس ویژگیهای آن اشیاء و ساختار پیوندهایشان انجام می‌دهد. یک حوزه کاربرد آن تعیین اجتماعات و بی<sup>۷</sup> است یعنی مجموعه‌ای از صفحات وب که بر یک عنوان یا زمینه<sup>۸</sup> خاص مرکز هستند. کاربرد مشابه، تعیین اجتماعات پژوهشی در حوزه کتاب‌شناسی می‌باشد.

**کشف زیر گرافها<sup>۹</sup>:** تعیین زیر گرافها، زیر گرافهای ویژه‌ای را در درون شبکه می‌یابد و نوعی از جستجوی گرافی است. مثالی از بیولوژی کشف زیر گرافهای متناظر با ساختار پروتئینهاست. در شیمی نیز می‌توان زیر گرافهایی را جستجو کرد که زیر ساختارهای شیمیایی را نمایش می‌دهند.

**فراداده‌کاوی:** فراداده‌ها، داده‌هایی در مورد داده‌ها هستند. فراداده‌ها، داده‌هایی نیمه‌ساخت‌یافته درباره داده‌های ساخت‌نیافته فراهم می‌کنند که طیف وسیعی از داده‌های وب و متن گرفته تا پایگاههای داده‌های چند رسانه‌ای را در بر می‌گیرد. این وظیفه برای یکپارچه‌سازی داده‌ها در

<sup>۱</sup>- Reconciliation

<sup>۲</sup>- Information Extraction

<sup>۳</sup>- Duplication

<sup>۴</sup>- Matching

<sup>۵</sup>- Identity Uncertainty

<sup>۶</sup>- Record Linkage

<sup>۷</sup>- Web Communities

<sup>۸</sup>- Theme

<sup>۹</sup>- Subgraph

بسیاری حوزه‌ها مفید است. فراداده کاوی را می‌توان برای نگاشت شماتیک<sup>۱</sup> (مثلاً ویژگی شماره مشتری از یک پایگاه داده به شماره مشتری از پایگاه داده‌ای دیگر نگاشت می‌شود چون هر دو آنها به یک موجودیت اشاره دارد؛ کشف شماتیک<sup>۲</sup> (از داده‌های نیمه‌ساخت یافته طرح‌هایی ایجاد می‌کند)؛ و تشکیل مجدد شماتیک<sup>۳</sup> (بر اساس فراداده‌های داده کاوی شده، طرح را اصلاح می‌کند) به کار گرفت. نمونه‌ها شامل این موارد است: انطباق دو منبع کتاب‌شناسخنی، کشف طرح از داده‌های نیمه‌ساخت یافته یا غیرساخت یافته روی وب، و نگاشت بین دو هستان‌شناسی<sup>۴</sup> دارویی.

به‌طور خلاصه، استفاده از اطلاعات پیوند بین اشیاء، وظیفه‌ای اضافی برای پیوند کاوی در مقایسه با رویکردهای سنتی کاوش به همراه می‌آورد. به کارگیری این وظایف در هر حال چالش‌های بسیاری را در برخواهد داشت. در اینجا چند نمونه از این چالش‌ها را بررسی می‌کیم.

**وابستگیهای منطقی در برابر وابستگیهای آماری:** دو نوع وابستگی در ساختارهای گراف قرار دارند: ساختارهای پیوند (میان رابطه منطقی بین اشیاء) و وابستگیهای احتمالی<sup>۵</sup> (میان رابطه آماری مثل همبستگی بین ویژگیهای اشیاء در حالتی که نوعاً چنین اشیایی منطقاً به هم مرتبطند). کار هم‌زمان با این وابستگیها نیز خود چالشی برای کاوش داده‌های چند رابطه‌ای است که در آن داده‌هایی که مورد کاوش قرار می‌گیرند، در جداول چند لایه‌ای قرار دارند. علاوه بر جستجوهای استاندارد بر روی وابستگیهای احتمالی بین ویژگیها، می‌بایست بر روی رابطه‌های منطقی مختلف ممکن بین اشیاء نیز جستجو صورت گیرد. این امر فضای جستجوی وسیعی را می‌طلبد که یافتن یک مدل ریاضی موجه را مشکل می‌سازد. روش‌های توسعه یافته در برنامه‌ریزی منطق قیاسی که بر روی ارتباطات منطقی جستجو می‌کند، اینجا می‌تواند مفید واقع شود.

**ساخت مشخصه‌ها:**<sup>۶</sup> در دسته‌بندی مبتنی بر پیوندها، ما هم به ویژگیهای یک شء توجه می‌کنیم و هم به ویژگیهای اشیایی متصل به آن. به علاوه، ممکن است پیوندها هم ویژگیهایی

<sup>۱</sup>- Schema Mapping

<sup>۲</sup>- Schema Discovery

<sup>۳</sup>- Reformulation Schema

<sup>۴</sup>- Ontologies

<sup>۵</sup>- Probabilistic Dependencies

<sup>۶</sup>- Feature Construction

داشته باشند. هدف از ساخت مشخصه‌ها، ایجاد یک صفت منفرد مبین این ویژگیهای است. این امر می‌تواند شامل انتخاب مشخصه‌ها<sup>۱</sup> و تجمعیت مشخصه‌ها<sup>۲</sup> باشد. در انتخاب مشخصه‌ها، فقط مشخصه‌های مهم و متمایزکننده درنظر گرفته می‌شوند. تجمعیت مشخصه‌ها، چند مجموعه از مقادیر بر روی مجموعه اشیاء مرتبط را می‌گیرد و خلاصه‌ای از آن را باز می‌گرداند. این خلاصه می‌تواند به عنوان مثال مُد مجموعه (مقداری که بیشترین تعداد رخداد را دارد)؛ میانگین مجموعه (اگر مقادیر عددی باشد)؛ یا میانه (اگر مقادیر به ترتیب مرتب شده باشند) باشد. برخی اوقات این روش مناسب نیست.

**نمونه‌ها<sup>۳</sup>** در برابر دسته‌ها: این چالش مربوط است به اینکه آیا این مدل صریحاً به تک تک نمونه‌ها اشاره دارد یا به دسته‌هایی (طبقه‌های عمومی) از نمونه‌ها. یک مزیت مدل پیشین در این است که می‌تواند برای اتصال تک تک نمونه‌های خاص با احتمال بالا به کار رود. یک مزیت مدل آخر این است که می‌تواند برای عمومیت بخشیدن به موقعیتها جدید با نمونه‌های منفرد مختلف استفاده شود.

دسته‌بندی جمعی و یکی‌سازی جمعی<sup>۴</sup>: آموزش دادن<sup>۵</sup> یک مدل برای دسته‌بندی را در نظر بگیرید که بر اساس مجموعه‌ای از اشیاء که برچسب یا نام دسته‌شان مشخص است، انجام شود. روش‌های دسته‌بندی سنتی تنها به ویژگیهای شیء توجه می‌نمود. فرض کنید پس از آموزش مدل، مجموعه جدیدی از اشیایی که برچسب دسته‌شان معلوم نیست در اختیار داریم. به کارگیری مدل برای تعیین برچسب دسته اشیاء جدید، با توجه به همبستگی‌های احتمالی بین اشیاء پیچیده است (برچسب دسته اشیاء به هم مرتبط ممکن است همبسته باشد). بنابراین دسته‌بندی می‌باشد گام تکرار شونده دیگری را هم در بر بگیرد که برچسب دسته هر شیء را بر اساس برچسب دسته اشیاء مرتبط با آن تغییر دهد (یا ثابت کند). در این معنا، دسته‌بندی بیشتر جمعی انجام می‌شود تا مستقل<sup>۶</sup>.

<sup>۱</sup>- Feature Selection

<sup>۲</sup>- Feature Aggregation

<sup>۳</sup>- Instances

<sup>۴</sup>- Collective Consolidation

<sup>۵</sup>- Training

**استفاده مؤثر از داده‌های برچسب خورده و برچسب نخورده:** یک استراتژی جدید در بادگیری این است که مخلوطی از داده‌های برچسب خورده و برچسب نخورده را شرکت دهد. داده‌های برچسب نخورده می‌تواند به استنتاج توزیع ویژگیهای اشیاء کمک کند. پیوند بین داده‌های برچسب خورده (داده‌های آموزشی) و برچسب نخورده (داده‌های آزمون) پی به وابستگیهایی می‌برد که می‌تواند به استدلالهای دقیق‌تر کمک کند.

**پیش‌بینی پیوندها:** یک چالش موجود در پیش‌بینی پیوندها این است که احتمال پیشین یک پیوند خاص بین اشیاء بسیار کم است. رویکردهای مختلف در باب پیش‌بینی پیوندها بر اساس تعدادی معیار برای تحلیل مجاورت گره‌ها در یک شبکه، مطرح شده‌اند. مدل‌های احتمالی هم مطرح شده‌اند. برای مجموعه‌های بزرگ داده ممکن است مدل کردن پیوندها در سطحی بالاتر مؤثر باشد

فرض دنیای باز در مقابل دنیای بسته<sup>۱</sup>: سنتی‌ترین رویکردها فرض می‌کنند که ما تمام موجودیتهای بالقوه را در این حوزه می‌شناسیم. این فرض دنیای بسته در کاربردهای دنیای واقع، غیرواقعی است. کار در این زمینه، معرفی یک زیان برای تعیین توزیعهای احتمال ساختارهای رابطه‌ای را دربرمی‌گیرد که متضمن مجموعه‌ای متغیر از اشیاء است.

**اجتماع‌کاوی بر روی شبکه‌های چندراطه‌ای:** کار بر روی تحلیل شبکه‌های اجتماعی نوعاً کشف گروههایی از اشیاء را که در ویژگیهای مشابهی سهیم هستند، دربرمی‌گیرد. به این کار اجتماع‌کاوی گفته می‌شود. پیوند<sup>۲</sup> صفحات وب یک مثال است که در آن اجتماع کشف شده می‌تواند مجموعه‌ای از صفحات وب مربوط به یک مبحث خاص باشد. اغلب الگوریتمهای اجتماع‌کاوی فرض می‌کنند که تنها یک شبکه اجتماعی وجود دارد که نشانگر رابطه نسبتاً همگنی است. در واقعیت، چندین شبکه اجتماعی ناهمگن وجود دارد که بین‌کار روابط گوناگونی هستند. چالش جدید کاوش اجتماعات پنهان در چنین شبکه‌های اجتماعی ناهمگن است که به

<sup>۱</sup>- Closed Versus Open World Assumption

<sup>۲</sup>- Community

<sup>۳</sup>- Linkage

آن اجتماع‌کاوی بر روی شبکه‌های اجتماعی چندرابطه‌ای گفته می‌شود. این چالشها ادامه دارد تا انگلیزهای برای تحقیقات بیشتر در پیوند کاوی باشد.

## ۴-۹- کاوش شبکه‌های اجتماعی

در این قسمت ما چند مثال در حوزه‌های مختلف کاوش بر روی شبکه‌های اجتماعی را بررسی می‌کنیم. این نمونه‌ها شامل پیش‌بینی پیوندها، کاوش شبکه‌های مشتریان برای بازاریابی ویروسی<sup>۱</sup>، کاوش گروه‌های خبری با استفاده از شبکه‌ها و اجتماع‌کاوی از شبکه‌های چندرابطه‌ای است. سایر نمونه‌ها شامل کشف زیرگرافهای مشخصه در گراف کاوی و کاوش ساختارهای پیوند در وب کاوی است. کاربردهای دیگری مانند خوشبندی و دسته‌بندی مبتنی بر پیوند نیز وجود دارند.

### پیش‌بینی پیوندها: چه یالهایی به شبکه افزوده خواهد شد؟

شبکه‌های اجتماعی پویا هستند. پیوندهای جدیدی ظاهر می‌شده که نشان دهنده تعاملی جدید بین اشیاء هستند. در مسئله پیش‌بینی پیوندها، یک تصویر آنی<sup>۲</sup> از شبکه اجتماعی در زمان  $t+1$  در اختیار ما قرار داده می‌شود و پیش‌بینی اینکه در بازه زمانی  $t$  تا  $t+1$  چه یالهایی به این شبکه افزوده خواهد شد، از ما خواسته می‌شود. در این حالت ما به دنبال این هستیم که با استفاده از صفات حقیقی خود مدل، از توسعه‌ای که می‌تواند تکامل یک شبکه اجتماعی را مدل کند، پرده برداریم. به عنوان مثال یک شبکه اجتماعی همنویسنده‌ی یا تألفی مشترک، بین دانشمندان را در نظر بگیرند. به طور شهودی ممکن است پیش‌بینی کنیم که دو دانشمند که در شبکه نزدیک<sup>۳</sup> یکدیگر قرار دارند، احتمال دارد در آینده با هم همکاری داشته باشند. بر این اساس پیش‌بینی پیوندها را می‌توان به عنوان بخشی از مطالعات مدل‌های تکامل شبکه‌های اجتماعی در نظر گرفت.

رویکردهایی که در باب پیش‌بینی پیوندها مطرح شده‌اند، مبتنی بر معیارهای مختلفی از تحلیل مجاورت گره‌های یک شبکه می‌باشند. بسیاری از معیارها از روش‌های تحلیل شبکه‌های

<sup>۱</sup>- Viral Marketing

<sup>۲</sup>- Snapshot

اجتماعی و نظریه گراف سرچشمه می‌گیرند. روش کلی بدین صورت است: در تمام روشها به هر جفت گره  $X$  و  $Y$ , بر اساس گراف وروودی  $G$  و اندازه مجاورت موجود، یک وزن اتصال<sup>۱</sup>  $Scroe(X, Y)$  تخصیص می‌دهند. سپس یک ترتیب نزولی از  $Scroe(X, Y)$  تهیه می‌شود که پیوندهای جدید پیش‌بینی شده را به ترتیب نزولی اطمینان به ما می‌دهد. این پیش‌بینیها را می‌توان بر اساس مشاهدات واقعی از مجموعه داده‌های تجربی ارزیابی نمود. ساده‌ترین رویکرد، جفتهای  $(X, Y)$  را براساس طول کوتاهترین مسیرشان در  $G$  مرتب می‌کند. این رویکرد نظریه دنیاهای کوچک را مجسم می‌کند که در آن تک تک اجزا از طریق کوتاهترین زنجیره به یکدیگر متصل هستند. از آنجا که هدف مشترک همه روشها، مرتب کردن تمام جفتهای به ترتیب کاهش امتیاز می‌باشد، در اینجا  $Scroe(X, Y)$  به صورت منفی طول کوتاهترین مسیر تعریف می‌شود. بسیاری از معیارها از اطلاعات همسایگی استفاده می‌کنند. ساده‌ترین نوع چنین معیارهایی همسایگان مشترک<sup>۲</sup> است (هر قدر تعداد همسایگان مشترک  $X, Y$  بیشتر باشد، احتمال اینکه در آینده بین  $X, Y$  پیوندی ایجاد شود بیشتر است). از لحاظ شهودی، اگر نویسنده  $X, Y$  هرگز مقاله مشترکی تالیف نکرده باشند ولی همکاران مشترک بسیاری داشته باشند، با احتمال بیشتری آنها در آینده با یکدیگر همکاری خواهند داشت. سایر معیارها بر اساس مجموع تمام مسیرهای<sup>۳</sup> بین دو گره استوار است. به عنوان مثال معیار کتز<sup>۴</sup> تمام مسیرهای وزن دهنده بین  $X$  و  $Y$  را محاسبه می‌کند به طوری که به مسیرهای کوتاهتر وزن بیشتری اختصاص می‌دهد. تمام این معیارها را می‌توان با ترکیبی از رویکردهای سطح بالاتر مانند خوشبندی به کار برد. به طور نمونه، روش پیش‌بینی پیوند را می‌توان برای نسخه اصلاح شده یک گراف به کار برد که در آن یالهای قلابی حذف شده‌اند.

در آزمایشات انجام شده بر روی مجموعه داده‌ها استناد یا نقل قول، هیچ‌کدام از روشها بر دیگر روشها مقدم نیست. بسیاری روشها به طور عمده یک پیشگویی تصادفی<sup>۵</sup> را به دست

<sup>۱</sup>- Connection Weight

<sup>۲</sup>- Common Neighbors

<sup>۳</sup>- Ensemble of All Paths

<sup>۴</sup>- Katz

<sup>۵</sup>- Random Predictor

می‌آورند که معتقد است توپولوژی شبکه‌ها می‌تواند اطلاعات مفیدی برای پیش‌بینی پیوندها فراهم آورد. معیار کتز و نسخه‌ها گوناگون مبتنی بر خوشبندی آن، همواره خوب عمل کرده‌اند هرچند دقیق پیش‌بینی همچنان بسیار پایین است. کارهای آینده در زمینه پیش‌بینی پیوندها ممکن است هم بر یافتن راههای بهتر استفاده از اطلاعات توپولوژی شبکه مرکز شود و هم کارآیی محاسبات فاصله گره‌ها را مثلاً از طریق تخمین زدن بهبود بخشد.

### کاوش شبکه‌های مشتریان به منظور بازاریابی ویروسی

بازاریابی ویروسی کاربردی از کاوش شبکه‌های اجتماعی است که مشخص می‌کند چگونه افراد می‌توانند بر رفتار خرید سایرین تأثیر بگذارند. به طور سنتی شرکتها، بازاریابی مستقیم<sup>۱</sup> (که در آن تصمیم فروش از طریق یک فرد خاص و صرفاً بر اساس خصوصیات آن فرد گرفته می‌شود) یا بازاریابی انبوه<sup>۲</sup> (که در آن افراد بر اساس اینکه به کدام بخش<sup>۳</sup> از جمعیت متعلق باشند، صورت می‌گیرد) را به کار برده‌اند. در هر حال این رویکردها تأثیری را که مشتریان می‌توانند بر تصمیم خرید سایرین داشته باشند نادیده می‌گیرند. به عنوان مثال فردی را در نظر بگیرید که تصمیم می‌گیرد یک فیلم خاص را ببیند و گروهی از دوستان را نیز برای تماشای آن فیلم تشویق می‌کند. هدف بازاریابی ویروسی بهینه کردن تأثیر مثبت گفتار شفاهی<sup>۴</sup> در میان مشتریان است. ممکن است بخواهیم هزینه بیشتری برای جذب یک نفر که تماسهای اجتماعی زیادی دارد، اختصاص دهیم. بنابراین با در نظر گرفتن این تعاملات بین مشتریان، بازاریابی ویروسی می‌تواند سود بیشتری نسبت به بازاریابی سنتی که چنین تعاملاتی را نادیده می‌گرفت، کسب نماید. رشد اینترنت در دهه‌های اخیر موجب در دسترس قرار گرفتن شبکه‌های اجتماعی بسیاری شده است که می‌توان با هدف بازاریابی ویروسی به کاوش آنها پرداخت، مثل لیستهای پست الکترونیک، گروه‌های خبری<sup>۵</sup>، اجتماعات برخط<sup>۶</sup> گپهای IRC<sup>۷</sup>، پیام آنی<sup>۸</sup>، سیستمهای فیلتر

<sup>۱</sup>- Direct Marketing

<sup>۲</sup>- Mass Marketing

<sup>۳</sup>- Segment

<sup>۴</sup>- Word of Mouth

<sup>۵</sup>- Use Net Groups

<sup>۶</sup>- On Line Forums

<sup>۷</sup>- Instant Relay Chat

مشارکتی و سایتهای شرکت دانش<sup>۲</sup>. سایتهای شرکت دانش (مانند *Epinions* در [www.epinions.com](http://www.epinions.com)) به کاربران خود اجازه می‌دهند (نوعاً به صورت مجانی) محصولات را به سایرین توصیه کنند یا برای محصولات برآورد قیمت نمایند تا به سایرین کمک کرده باشند. کاربران می‌توانند مفید بودن یا قابلیت اعتماد<sup>۳</sup> یک مقاله مروری<sup>۴</sup> را ارزش‌گذاری کنند، و حتی در صورت امکان سایر متقدان را نیز ارزیابی نمایند. در این صورت یک شبکه روابط با عنوان شبکه اعتماد<sup>۵</sup> بر اساس اعتماد بین کاربران شکل می‌گیرد که بیانگر یک شبکه اجتماعی است که می‌تواند مورد کاوش قرار بگیرد.

ارزش شبکه‌ای<sup>۶</sup> یک مشتری، افزایش مورد انتظار در فروش به سایرین است که از جذب آن مشتری حاصل می‌شود. در مثال فوق، اگر مشتری مورد نظر ما سایرین را متقاعد کند که یک فیلم خاص را ببینند، شرکت فیلم‌سازی به صرف هزینه بیشتر برای ترغیب وی به مشاهده فیلم مشتاق می‌شود. در عوض اگر این مشتری نوعاً هنگام تصمیم‌گیری درباره اینکه چه فیلمی را تماشا کند، به سایرین گوش می‌کند، هزینه بازاریابی صرف شده برای وی را می‌توان اتلاف منابع محسوب نمود. بازاریابی ویروسی ارزش شبکه‌ای یک مشتری را در نظر می‌گیرد. در حالت ایده‌آل، ترجیح می‌دهیم شبکه یک مشتری را مورد کاوش قرار دهیم (به عنوان مثال شبکه دوستان و اقوام وی) تا نه تنها بر اساس ویژگیهای مشتری، بلکه بر اساس تأثیر همسایگان مشتری در شبکه، پیش‌بینی کنیم، چقدر احتمال دارد وی یک کالای خاص را بخرد. اگر مجموعه مشخصی از مشتریان را مورد بازاریابی قرار دهیم، از طریق بازاریابی ویروسی می‌توانیم میزان سود مورد انتظار از کل شبکه را پس از آنکه تأثیر آن مشتریان در شبکه تکثیر شد، در تحقیقاتمان دنبال کنیم. این کار می‌تواند به ما در یافتن مجموعه بهینه مشتریانی که مورد بازاریابی قرار می‌گیرند، یاری رساند. ارزش شبکه‌ای مشتریان (که در بازاریابی مستقیم سنتی نادیده گرفته می‌شوند) می‌تواند به طرح بازاریابی بهبود یافته‌ای بیانجامد.

<sup>۱</sup>- Instant Messaging

<sup>۲</sup>- Knowledge- Sharing Sites

<sup>۳</sup>- Trust Worthiness

<sup>۴</sup>- Review

<sup>۵</sup>- Web of Trust

<sup>۶</sup>- Network Values

مجموعه  $n$  مشتری بالقوه را در نظر بگیرید،  $X_i$  را یک متغیر دو حالته فرض کنید که اگر مشتری  $i$  محصول بازاریابی شده را بخرد مقدار یک می‌گیرد و در غیر این صورت مقدارش صفر می‌شود. همسایگان  $X_i$  مشتریانی هستند که مستقیماً بر  $X_i$  تأثیر می‌گذارند.  $M_i$  اقدام بازاریابی<sup>۱</sup> است که در مورد مشتری  $i$ م صورت می‌گیرد.  $M_i$  می‌تواند دو حالته (اگر کوپنی<sup>۲</sup> برای این مشتری ارسال شود، یک و در غیر این صورت صفر) یا طبقه‌ای (نشانگر اقدامات مختلف ممکن اتخاذ شده در قبال مشتری) باشد. در حالتی دیگر  $M_i$  می‌تواند مقداری پیوسته (مثلًا نشان دهنده میزان تخفیفی که برای مشتری قائل شده‌ایم) باشد. یافتن طرح بازاریابی که سود را بیشینه کند، مطلوب ماست. یک مدل احتمالی مطرح شده است که  $M_i$  را به عنوان مقداری پیوسته، بهینه می‌کند. در این مدل به جای اینکه تنها تصمیمی دو حالته اتخاذ شود که آیا این مشتری را مورد بازاریابی قرار بدهد یا خیر، مقدار هزینه صرف شده برای بازاریابی هر مشتری، بهینه می‌شود.

این مدل عواملی را که بر ارزش شبکه مؤثرند، در نظر می‌گیرد. اول این مشتری می‌بایست اتصالات زیادی در شبکه داشته باشد و به علاوه محصول در نظر وی خوب قلمداد شود. در صورتی که یک مشتری دارای اتصالات زیاد، نگرش منفی نسبت به محصول مورد نظر داشته باشد، ارزش شبکه‌ای او می‌تواند منفی باشد که در این حالت، بازاریابی (جذب) وی توصیه نمی‌شود. دوم، این مشتری می‌باید ترجیح‌آور دیگران تأثیرگذار باشد تا تأثیرپذیر. سوم طبیعت بازگشتی<sup>۳</sup> اثر گذاری تبلیغ شفاهی را می‌بایست مد نظر قرار داد. یک مشتری ممکن است بر آشنایان خود تأثیر بگذارد. آنها هم به نوبه خود ممکن است محصول را پیشنهاد و بر سایر افراد تأثیر بگذارند وضع به همین ترتیب ادامه می‌باید تا جایی که کل شبکه پوشش داده می‌شود. در ضمن این مدل ملاحظات مهم دیگری را نیز در نظر می‌گیرد: ممکن است تولید کننده بپذیرد که مقداری پول را برای مشتریانی که به اندازه کافی بر دیگران تأثیر مثبت دارند از دست بدهد. به عنوان مثال، فروش مجاني به یک مشتری برگزیده مناسب می‌تواند چندین برابر در فروش به سایر مشتریان جبران شود. این رویکرد چرخشی بزرگ نسبت به بازاریابی سنتی مستقیم است که در آن در صورتی که سود حاصل از فروش به یک مشتری به تنها بیان از هزینه توصیه محصول به

<sup>۱</sup>- Marketing Action

<sup>۲</sup>- Coupon

<sup>۳</sup>- Recursive

وی تجاوز کند، تحفیفی به آن مشتری تعلق می‌گیرد. این حقیقت را که، تنها دانش جزئی نسبت به شبکه داریم و جمع‌آوری چنین دانشی نیز می‌تواند هزینه‌های مربوط به خود را داشته باشد، در این مدل مورد توجه قرار می‌گیرد.

یافتن مجموعه مشتریان بهینه به صورت یک مسئله بهینه‌سازی که به خوبی تعریف شده، فرموله شده است: یافتن مجموعه مشتریانی که سود خالص را بیشینه می‌کنند. این مسئله از نوع دشوار  $NP$  شناخته شده است. با به کارگیری رویه جستجوی تپه‌نوردي ساده می‌توان با تقریب ۶۳٪ جواب بهینه را به دست آورد. مادامی که افزایش مشتریان سود سرانه را بهبود بخشد، مشتریان جدیدی به مجموعه مشتریان افزوده خواهند شد. این روش با وجود دانش ناقص از شبکه باثبات<sup>۱</sup> شناخته شده است.

روشهای بازاریابی ویروسی را می‌توان برای حوزه‌های دیگر نیز به کار برد. کاهش انتشار ویروس  $HIV$  مبارزه با استعمال دخانیات در نوجوانان و ابتکار سیاسی اجتماعی محلی<sup>۲</sup> مثالهایی از این نوع هستند. به کارگیری روشهای بازاریابی ویروسی برای حوزه وب و عکس آن، زمینه‌های جالبی برای تحقیقات آنی هستند.

#### ۱-۴-۹ - کاوش گروه‌های خبری با کمک شبکه‌ها

تحلیل شبکه‌های اجتماعی مبتنی بر وب با وب‌کاوی رابطه نزدیکی دارد. در وب‌کاوی دو الگوریتم رتبه‌بندی متداول به نام صفحه‌رتبه و  $HITS$  به کار می‌رود. این الگوریتمها بر این پایه استوارند که پیوندی از صفحه وی  $A$  به  $B$  معمولاً نشانه تأیید  $B$  توسط  $A$  است.

وضعیت در گروه‌های خبری روی مباحث موضوعی متفاوت است. یک پست نوعی در گروه خبری شامل یک یا چند خط نقل قول<sup>۳</sup> از پستی دیگر و به دنبال آن نظر نویسنده پست فعلی است. این پاسخ‌های نقل قول دار پیوندهای نقل قولی را شکل داده و شبکه‌ای را ایجاد می‌کنند که در آن رئوس، بیانگر افراد و پیوندها رابطه پاسخ هستند. پدیده جالب این است که مردم بیشتر به پیامی که با آن مخالفند پاسخ می‌دهند تا پیامی که موافقند. این رفتار که در بسیاری از گروه‌های

<sup>۱</sup>- Robust

<sup>۲</sup>- Grass- Root Political Initiative

<sup>۳</sup>- Quote

خبری وجود دارد کاملاً با گراف پیوند صفحات وب که در آن پیوند نشانه توافق یا علاقه مشترک می‌باشد، متضاد است. بر پایه این رفتار می‌توان با تحلیل ساختار گراف پاسخ‌ها به طور مؤثر، نویسنده‌گان داخل گروه خبری را به دسته‌های مخالف دسته‌بندی و افزایش کرد. این فرایند دسته‌بندی گروه خبری، با نظریه گراف قابل انجام است. اگر فرد  $\alpha$  از پست قبلی فرد  $\beta$  نقل قول کند، پیوند نقل قول بین  $\alpha$  و ز ساخته شده و از روی این پیوندها شبکه یا گراف نقل قول ایجاد می‌شود. حال دوبخشی کردن رئوس به دو مجموعه را در نظر می‌گیریم:

مجموعه  $F$  بیانگر موافقین یک مطلب و مجموعه  $A$  بیانگر مخالفین آن مطلب است. اگر يالهای گراف گروه خبری بیانگر مخالفت باشند آن گاه انتخاب بهینه، حداکثر کردن تعداد يالها بین این دو مجموعه است. از آنجا که مسئله حداکثر برش (حداکثر کردن تعداد يالها برش خورده برای دو بخش کردن گراف) به طور نظری مسئله‌ای از نوع  $NP$  است، نیاز به راه حل عملی دیگری داریم. به خصوص می‌توان از دو حقیقت دیگر در وضعیت فعلی استفاده کرد: نمونه ما بیشتر از آن که گراف عامی باشد، گرافی دوبخشی است که برخی گرهای مغشوش به آن اضافه شده‌اند.

هیچ‌کدام از دو بخش چندان از دیگری کوچکتر نیستند.

در چنین وضعیتهايی می‌توان مسئله را به مسئله برش تقریباً متوازن حداقل وزن تبدیل کرد، که به نوبه خود با روش‌های طیفی ساده تقریب زده می‌شود. برای بهبود دقت دسته‌بندی می‌توان ابتدا به طور دستی تعداد کمی از پست‌کننده‌گان فعال را دسته‌بندی کرد و رئوس متناظر در گراف را علامت زد. سپس از این اطلاعات برای دستیابی به افزایی بهتر استفاده کرد به این نحو که در حین اجرای الگوریتم افزایی، قرار داشتن موارد دستی در دو سو حفظ شود.

بر پایه این ایده‌ها، الگوریتم مؤثری ارائه شده است. آزمایشاتی با چند مجموعه داده گروه خبری در مباحث اجتماعی بحث‌انگیز مانند سقط جنین، کنترل اسلحة و مهاجرت نشان می‌دهد که پیوندها حاوی اطلاعاتی کم‌اغتشاش‌تر از متون هستند. روش‌های مبتنی بر تحلیل زبانی و آماری متن، صحت کمتری از تحلیل پیوند در این نوع گروه‌های خبری داشتند، زیرا الفاظ مورد استفاده طرفهای مقابله تا حد زیادی مشابه بوده و بسیاری از پستها حاوی متنی بسیار مختصر هستند که امکان تحلیل زبانی قابل اعتمادی نمی‌دهد.

## ۲-۴-۹- اجتماع کاوی شبکه‌های چندرابطه‌ای

با رشد وب، اجتماع کاوی توجهات روز افزونی را به خود جلب نموده است. قسمت اعظم چنین کارهایی بر کاوش اجتماعات ضمنی صفحات وب، اجتماعات ضمنی ادبیات علمی برگرفته از وب و اجتماعات ضمنی استنادات متمرکز شده است. در اصل یک اجتماع را می‌توان به صورت گروهی از اشیاء که در چندین خصوصیات مشترک سهیم هستند، تعریف نمود. اجتماع کاوی را می‌توان به صورت تعیین زیر گرافها در نظر گرفت. به عنوان مثال، در صفحات وب مربوط، دو صفحه وب (اشیاء) در صورتی به هم مرتبط هستند که ابرپیوندی بین آنها وجود داشته باشد. گرافی از روابط صفحات وب را می‌توان به منظور تعیین اجتماع یا تعیین مجموعه‌ای از صفحات وب در مورد یک موضوع خاص، مورد کاوش قرار داد.

اغلب روش‌های گراف‌کاوی و اجتماع کاوی مبتنی بر گرافهای همگن است. یعنی آنها فرض می‌کنند تنها یک نوع رابطه بین اشیاء وجود دارد. در حالی که در شبکه‌های اجتماعی واقعی، همواره انواع مختلف روابط بین اشیاء برقرار است. به هر نوع رابطه می‌توان در قالب یک شبکه روابط<sup>۱</sup> نگریست. (که به آن شبکه‌های اجتماعی همگن نیز گفته می‌شود). در این صورت چند نوع رابطه یک شبکه اجتماعی چندرابطه‌ای<sup>۲</sup> را تشکیل می‌دهند (که به آن شبکه‌های اجتماعی غیر همگن نیز گفته می‌شود). هر نوع رابطه ممکن است نقش معینی در یک وظیفه خاص ایفا نمایند. در اینجا گرافهایی با روابط مختلف می‌توانند برای ما اجتماعات متفاوتی را ایجاد نمایند.

به منظور یافتن اجتماعی با خصوصیات معلوم ابتدا نیاز است تعیین شود، کدام رابطه نقشی مهم در چنین اجتماعی ایفا می‌کند. چنین رابطه‌ای ممکن است به طور صریح وجود نداشته باشد، یعنی ما نیاز داشته باشیم قبل از آنکه آن اجتماع را در چنان شبکه رابطه‌ای بیاییم ابتدا چنین رابطه پنهانی را کشف کنیم. کاربران مختلف ممکن است به روابط متفاوتی در درون شبکه علاقه‌مند باشند. بدین لحاظ اگر ما شبکه‌ها را تنها با در نظر گرفتن تنها یک نوع رابطه مورد کاوش قرار دهیم، ممکن است اطلاعات ارزشمند بسیاری از اجتماع پنهان را از دست بدهیم. چنین کاوشی نمی‌تواند نیازهای متنوع اطلاعاتی کاربران مختلف را برآورده کند. این امر ما را به مسئله کاوش

<sup>۱</sup>- Relation Network

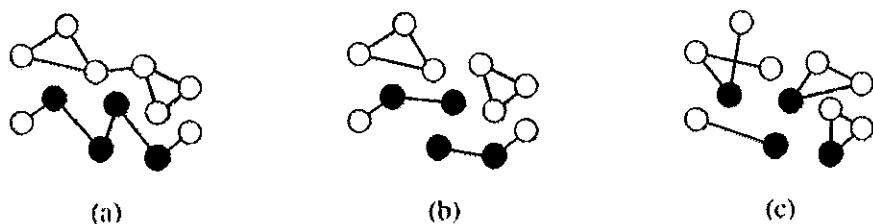
<sup>۲</sup>- Multi Relational Social Network

اجتماعات چندرابطه‌ای می‌رساند که کاوش اجتماعات پنهان در شبکه‌های اجتماعی ناهمگن را در بر می‌گیرد.

مثالی ساده را بررسی می‌کنیم. در یک اجتماع انسانی، ممکن است چندین نوع رابطه وجود داشته باشد: برخی از مردم در یک محل کار می‌کنند؛ بعضی علایق مشترکی دارند؛ برخی به یک درمانگاه می‌روند، والی آخر، از لحاظ ریاضی این اجتماع را می‌توان با یک گراف بزرگ نمایش داد که در آن گره‌ها نماد افراد هستند و يالها قوت رابطه بین آنها را مشخص می‌کنند. از آنجا که انواع مختلف رابطه وجود دارد، يالهای این گراف می‌بایست غیر همگن باشند. در برخی موارد، می‌توان این اجتماع را با به کار گیری چندین گراف همگن نیز مدلسازی نمود. هر گراف یک نوع رابطه را منعکس می‌کند. فرض کنید یک بیماری مسری شیوع پیدا کند، و دولت تلاش می‌کند افرادی را که بیشترین احتمال بیمار شدن را دارند، شناسایی نماید. به وضوح پیداست که روابط موجود بین مردم نمی‌توانند نقش یکسانی داشته باشند. منطقی است که فرض کنیم در چنین شرایطی رابطه «کارکردن در یک محل» یا «زندگی کردن با یکدیگر» می‌بایست نقشی حیاتی اینها نمایند. سوالی که پیش می‌آید این است: «چگونه می‌توانیم رابطه‌ای که بیشترین نقش را در شیوع بیماری دارد، انتخاب کنیم؟ آیا رابطه پنهانی (بر اساس روابط صریح و آشکار) وجود دارد که به بهترین نحو مسیر شیوع بیماری را فاش نماید؟»

سؤال از لحاظ ریاضی می‌تواند به صورت انتخاب و استخراج رابطه<sup>۱</sup> در تحلیل شبکه‌های اجتماعی چندرابطه‌ای، مدلسازی شود. مسئله استخراج رابطه را می‌توان به صورت ساده آن‌گونه که در ادامه خواهد آمد، بیان نمود: در یک شبکه اجتماعی چندرابطه‌ای، بر اساس چند نمونه برچسب گذاری شده است (به عنوان مثال به صورت پرس‌وجوهایی که کاربر فراهم می‌کند). چگونه می‌توان اهمیت روابط مختلف را ارزیابی نمود؟ در ضمن چگونه می‌توان به ترکیبی از روابط موجود دست یافت که بیشترین تطابق را با رابطه نمونه‌های برچسب خورده داشته باشد؟ به عنوان مثال شبکه نشان داده شده در شکل (۹-۳) را ببینید که سه نوع رابطه مختلف دارد؛ به ترتیب با (a)، (b) و (c) مشخص شده‌اند. فرض کنید کاربری نیاز دارد که ۴ شیء رنگی به یک اجتماع تعلق داشته باشند و این امر را با یک پرس‌وجو مشخص می‌کند. پیداست که اهمیت

نسبی هر کدام از سه نوع رابطه با توجه به نیاز اطلاعاتی کاربر فرق می‌کند. در بین این سه نوع رابطه، (a) بیشترین تطبیق را با نیاز کاربر دارد و لذا مهم‌ترین است، در حالی که (b) در رتبه دوم قرار می‌گیرد. رابطه (c) را می‌توان با توجه به نیاز اطلاعاتی کاربر، یک اختشاش در نظر گرفت. تحلیل سنتی شبکه‌های اجتماعی بین این روابط تمایزی قائل نمی‌شود و با روابط مختلف یکسان برخورد می‌کند. آنها به سادگی برای توصیف ساختار بین اشیا با یکدیگر ترکیب می‌شوند. متأسفانه، در این مثال، رابطه (c) تأثیری منفی بر این هدف می‌گذارد. با این وجود، اگر ما این روابط را بر اساس اهمیتشان ترکیب کنیم، رابطه (c) به آسانی حذف می‌شود و روابط (a) و (b) را برای کشف ساختار اجتماع باقی می‌گذارد که با نیاز کاربر مطابق است.



شکل ۳-۹ شبکه‌های با روابط مختلف

بعضی مواقع یک کاربر ممکن است پرس‌وجوی پیچیده‌تری را ارائه دهد. مثلاً ممکن است مشخص کند که دو شیء رنگی پایینی می‌باشد به دو اجتماع متفاوت متعلق باشند. در این صورت، اهمیت سه نوع رابطه شکل (۳-۹) تغییر می‌کند. رابطه (b) مهم‌ترین رابطه می‌شود، درحالی که رابطه (a) بی فایده قلمداد می‌شود (حتی با توجه به پرس‌وجو اثری منفی دارد). درنتیجه، در شبکه‌های اجتماعی چندرابطه‌ای، اجتماع کالوی می‌باشد بر اساس پرس‌وجوی کاربر (یا اطلاعاتی که لازم دارد) صورت بگیرد. پرس‌وجوی یک کاربر می‌تواند بسیار منعطف باشد. روش‌های اولیه تنها بر یک شبکه رابطه مرکز بودند و بر اساس پرس‌وجوی کاربر انجام نمی‌شدند و لذا نمی‌توانستند پاسخگوی چنین موقعیتهای پیچیده‌ای باشند.

الگوریتمی برای استخراج و انتخاب رابطه مطرح شده که مسئله را به صورت یک مسئله بهینه‌سازی مدل می‌کند. این مسئله از لحاظ ریاضی می‌تواند بدین صورت تعریف شود: مجموعه‌ای از اشیاء و مجموعه‌ای از روابط را در اختیار داریم که به صورت مجموعه‌ای از

گرافهای  $G_i(V, E_i)$ ،  $i = 1, \dots, n$  در آن  $n$  تعداد روابط است،  $V$  مجموعه گره‌های (اشیاء) و  $E_i$  مجموعه یالهای مربوط به نامین رابطه است. وزن یالها را می‌توان به صورت طبیعی بر اساس قوت رابطه بین دو شیء تعیین نمود. این الگوریتم هر رابطه را با یک گراف و یک ماتریس اوزان مشخص می‌کند.

$M_i$  را نماد ماتریس اوزان مربوط به  $G_i$  در نظر بگیرید. هر مؤلفه در ماتریس قوت رابطه بین یک جفت شیء مربوط را مشخص می‌کند. فرض کنید یک رابطه پنهانی با گراف  $(G^{\wedge}, E^{\wedge})$  مشخص شود و  $M^{\wedge}$  هم ماتریس اوزان مرتبط با  $G^{\wedge}$  است.

یک کاربر نیاز اطلاعاتی خود را به صورت یک پرس‌وجو مشخص می‌کند که در آن مجموعه‌ای از اشیاء برچسب‌گذاری شده  $[x_1, \dots, x_m] = X$  را اعلان می‌نماید.

$[y_1, \dots, y_m] = Y$  نیز به گونه‌ای است که در آن را برچسب  $x_i$  است (چنین اشیاء برچسب‌گذاری شده‌ای اطلاعات جزئی از روابط پنهان  $G^{\wedge}$  را نشان می‌دهند). هدف این الگوریتم یافتن ترکیب خطی این ماتریسهای اوزان است، به نحوی که به بهترین صورت  $G^{\wedge}$  (ماتریس اوزان مربوط به اشیاء دارای برچسب) را تخمین بزنند. ترکیب حاصله با احتمال بیشتری نیازمندی اطلاعاتی کاربر را برآورده ساخته و لذا موجب عملکرد بهتری در اجتماع کاوی خواهد شد.

این الگوریتم روی داده‌های کتاب‌شناسی آزمایش شده است. به طور طبیعی روابط چندگانه‌ای بین نویسنده‌گان وجود دارد. نویسنده‌گان می‌توانند مقالات خود را در هزاران کنفرانس مختلف به چاپ برسانند، و هر کنفرانس را می‌توان به صورت یک رابطه در نظر گرفت که یک شبکه اجتماعی چندرابطه‌ای را ایجاد خواهد کرد. با در دست داشتن چند مثال به دست آمده از کاربران (مانند گروه‌های نویسنده‌گان)، این الگوریتم می‌تواند یک رابطه جدید را با استفاده از مثال‌ها استخراج کند و تمام گروه‌های مرتبط دیگر را بیابد. رابطه استخراجی می‌تواند به صورت گروهی از نویسنده‌گان باشد که علایق مشترکی را دنبال می‌کنند.

## مراجع

- 1) Han, J, Kamber, M. (2006) "Chapter 9: Graph Mining, Social Network Analysis, and Multirelational Data mining", *Data mining concepts and techniques*, 2nd edition, , Morgan Kaufmann Publishers



---

## فصل دهم

---

# کاربرد داده‌کاوی در مدیریت ارتباط با مشتری

توجه: مطالب این فصل به طور مستقل از فصول قبل قابل مطالعه بوده و برای دانشجویان رشته‌های تجارت الکترونیک و مدیریت بازرگانی مناسب می‌باشد.

در طی چند سال گذشته تعامل شرکتها با مشتریانشان به طور قابل توجهی تغییر کرده است به طوری که تداوم کسب و کار با مشتری تضمین بلند مدت ندارد. به همین دلیل برای موفقیت یک سازمان لازم است سازمان‌ها مشتریانشان را به درستی درک کرده، نیازها و خواسته‌های آنها را پیش‌بینی کنند و با مجهز شدن به این اطلاعات، سلامت کاری خود را بهبود بخشنند. بسیاری از سازمانها داده‌های بسیار زیادی را درباره مشتریان، تامین کنندگان و شرکای تجاری‌شان جمع‌آوری و ذخیره می‌کنند ولی ناتوانی این سازمانها برای کشف داشن پنهان بالرزش در این داده‌ها سبب می‌شود که این داده‌ها به داشن تبدیل نشوند و این کار عملاً بیهوده باشد. صاحبان کسب و کار می‌بینند که استخراج اطلاعاتی ناشناخته، معتبر و قابل درک از بانک‌های اطلاعاتی عظیم خود و استفاده از این اطلاعات برای کسب سود بیشتر دارند. برای برخی، داده‌کاوی از نظر فنی جالب است ولی برای بیشتر مردم، این علم وسیله‌ای برای رسیدن به نتایج جالب می‌باشد. داده‌کاوی به تنها مفید نیست، بلکه زمانی که به صورت کاربردی در یک مورد خاص استفاده می‌شود، معنا پیدا می‌کند. برای محقق شدن این اهداف سازمانها باید مراحل زیر را طی نمایند:

- جمع‌آوری و یکپارچه‌سازی داده‌های داخلی و خارجی (خرید) در کل سازمان به شکلی قابل درک.
- کاوش داده‌های یکپارچه برای تولید دانش.
- سازماندهی و ارائه اطلاعات و دانش به شیوه‌ای که فرآیندهای تصمیم‌گیری پیچیده را تسريع نماید.

برای محقق شدن همه این اهداف، سازمانها نیاز به یکپارچه‌سازی مؤلفه‌های مختلف برنامه‌های کاربردی خود دارند. یکی از حوزه‌هایی که به سرعت در حال رشد است فناوری تصمیم‌گیری علمی<sup>۱</sup> (عمولاً به آنها موتورهای تصمیم‌گیری گفته می‌شود) و داده‌کاوی در مدیریت ارتباط با مشتری است.

برای حفظ رقابت، سازمانها نیاز به تدوین استراتژیهای تمرکز بر مشتری<sup>۲</sup>، مشتری محوری<sup>۳</sup>، مشتری مداری<sup>۴</sup> دارند. همه این موارد خواسته‌های سازمانها را در راستای ارتباط با مشتریان تعریف می‌کند. مدیریت ارتباط با مشتری<sup>۵</sup> راه حلی است که این تلاشها را برای سازمانها و همچنین مشتریان محقق می‌سازد.

فرض اولیه این است که همه مشتریان با یکدیگر برابر نیستند. هدف اصلی CRM، بهینه‌سازی نسبت وقایع سودآور به وقایع زیانبار برای گروه معلومی از مشتریان است. برخی وقایع خاص مانند فروش محصول، ایجاد درآمد و بعضی تولید سود می‌کنند در حالی که برخی دیگر مانند تماسهای تلفنی و دیگر موارد موارد مشابه به منظور پشتیبانی این‌گونه نیستند. در این‌گونه موارد هدف داده‌کاوی افزایش درآمد و کاهش هزینه می‌باشد. چیزی که باید سازمانها بدانند این است که مشتریان به چه شیوه‌ای و چگونه تمایل به تعامل دارند تا بتوانند وفاداری مشتری را کسب نموده و به طبع سودآوری خود را تبیز بهبد بخشنند. مدیریت ارتباط با مشتری به سازمانها این اجازه را می‌دهد که مشتریان خود را بهتر بشناسند و تفاوت بین آنها را بهتر درک

<sup>۱</sup>- Decision Science

<sup>۲</sup>- Customer Focused

<sup>۳</sup>- Customer Driven

<sup>۴</sup>- Customer Centric

<sup>۵</sup>- Customer Relationship Management (CRM)

نمایند. در نتیجه در تخصیص منابع به مشتریان مطلوب‌تر، کارآمدی بیشتری داشته باشد. از طریق تلاش‌های CRM، سازمانها می‌توانند هماهنگی بهتری در ارتباط با مشتری ایجاد نمایند، بنابراین سازمان می‌تواند مدیریت مؤثرتری بر روی منابع بازاریابی و ارتباط با معناتری با مشتریانش داشته باشد. ارتباط کارا با مشتری نیاز به درک الزامات این رابطه دارد. توانایی ارائه خدمات شخصی شده، ایجاد ارزش و پذیرش دوطرفه، تعهد به ارتباط متقابل و موارد مشابه همه در ایجاد ارتباط قوی تأثیر بهسزایی دارند.

## ۱۰-۱-معماری مدیریت ارتباط با مشتری

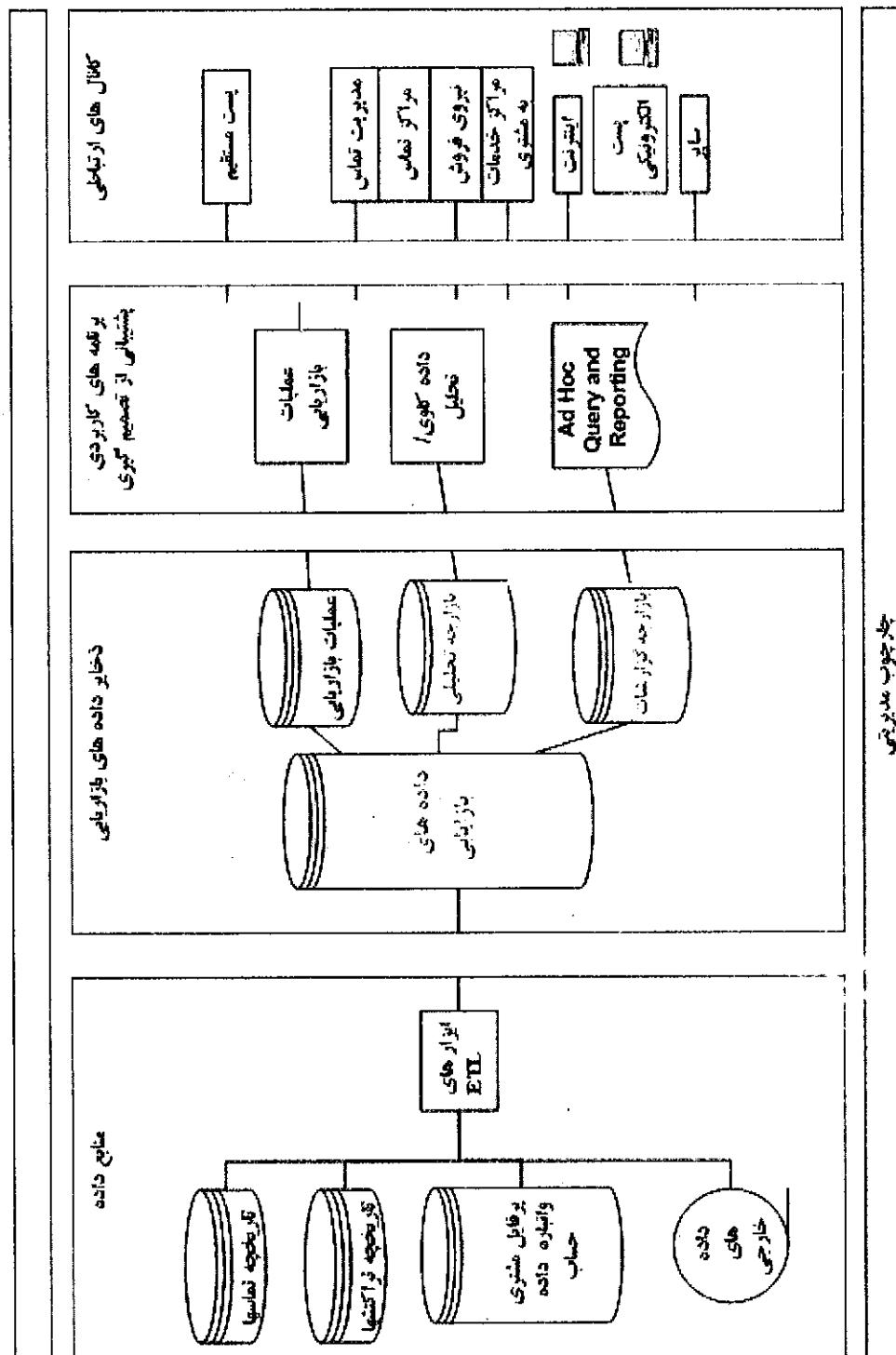
از دید معماری، کل چارچوب کاری CRM به سه مؤلفه اصلی تقسیم می‌شود:

- CRM عملیاتی<sup>۱</sup>: عبارتست از خودکارسازی فرآیندهای کسب‌وکاری افقی، شامل نقاط تماس مشتری، کانال‌ها و یکپارچه کردن موضوعات پشت صحنه و با موارد قابل مشاهده در روی صحنه.
  - CRM تحلیلی<sup>۲</sup>: تحلیل داده‌های ایجاد شده توسط CRM عملیاتی می‌باشد.
  - CRM مشارکتی<sup>۳</sup>: برنامه‌های کاربردی خدمات مشارکتی شامل موارد زیر می‌باشد: پست الکترونیکی، انتشارات شخصی شده، ارتباطات الکترونیکی، رسانه‌های مشابهی که برای تسهیل تعامل بین مشتری و سازمان طراحی شده است.
- همان‌طور که در (۱۰-۱) مشاهده می‌کنید معماری CRM شامل نقاط تماس با مشتری و کانال‌های ارسال می‌باشد که اطلاعات را تولید و مصرف می‌نمایند. این اطلاعات نیاز به یکپارچگی و تحلیل به منظور تعیین تصویر کامل و دقیق از مشتریان، عملکرد، نیازها، شکایات و مشخصه‌های آنها دارد [۳].

<sup>۱</sup>- Operational CRM

<sup>۲</sup>- Analytical CRM

<sup>۳</sup>- Collaborative CRM



نمکل ۱۰-۱) معماری CRM

### ۱-۱-۱- یافتن مشتریان احتمالی

هر یک از اهداف کسب و کار به فنون داده‌کاوی مناسب با آن مسئله مربوط می‌شوند. مباحث کسب و کار مورد بررسی مناسب با به ترتیب پیچیدگی رابطه مشتری می‌باشد. این بحث در ارتباط با مشتریان احتمالی<sup>۱</sup> که کمتر شناخته شده‌اند شروع شده و تا فرصت‌های متنوع بیان شده توسط روابط مستمر مشتری ادامه می‌یابد. موضوعات مورد نظر می‌تواند شامل محصولات، کانال‌های ارتباطی متعدد و تعاملات فردی فراینده باشد[۴].

به نظر می‌رسد جستجوی مشتریان احتمالی، شروع مناسبی برای بحث در مورد کاربردهای داده‌کاوی در این مبحث باشد. در بازاریابی، یک مشتری احتمالی کسی است که اگر به روش درستی با او برخورد شود، انتظار می‌رود به یک مشتری خوب تبدیل شود.

برای بیشتر کسب و کارها، تعداد کمی از شش میلیارد مردم زمین، مشتری احتمالی هستند. می‌توان بر اساس محل، سن، توانایی مالی و نیاز به محصول یا خدمت، اکثر مردم را کنار گذاشت. به عنوان مثال شرکتی که صندلی تاب حیاط می‌فروشد، قاعده‌تاً کاتالوگ خود را به خانوارهای دارای فرزند در مناطقی که حیاط خلوت دارند، ارسال می‌کند. یک مجله، گروهی را هدف‌گیری می‌کند که آن را مطالعه کرده و به آگهی‌هایش علاقه‌مند باشند.

داده‌کاوی می‌تواند نقش مهمی در یافتن مشتریان احتمالی داشته باشد. مهم‌ترین نقشهای آن عبارتند از:

- تشخیص مشتریان احتمالی خوب
  - انتخاب کanal ارتباطی مناسب به منظور دسترسی به مشتریان احتمالی
  - انتخاب پیام مناسب برای گروه‌های مختلف مشتریان احتمالی
- داده‌کاوی تا کنون بیشتر در تشخیص مشتریان احتمالی خوب استفاده شده است.

#### تشخیص مشتریان احتمالی خوب

بنابر ساده‌ترین تعریف، «مشتری احتمالی خوب» کسی است که ممکن است علاقه به مشتری شدن نشان دهد. مشتریان احتمالی خوب نه فقط علاقه‌مند هستند تا مشتری شوند، بلکه

استطاعت مالی مشتری شدن را نیز داشته و به عنوان مشتری سودآور خواهند بود. این مشتریان به احتمال زیاد صورتحساب‌هایشان را به موقع پرداخته، احتمال کمی دارد که سر شرکت کلاه بگذارند، اگر با آنها درست رفتار شود مشتریان وفاداری<sup>۱</sup> خواهند بود و دیگران را نیز همراه خواهند کرد. جدا از سادگی و پیچیدگی تعریف مشتری احتمالی خوب، اولین وظیفه هدف‌گیری آنها است.

اگر قرار است پیام از طریق تبلیغ یا کانالهای مستقیم‌تر مانند پست، تلفن یا پست الکترونیک انتقال یابد، هدف‌گیری<sup>۲</sup> بسیار مهم است. پیامهای تابلوهای تبلیغاتی تا حدودی هدف‌گیری شده‌اند. مثلاً تابلوهای تبلیغاتی خطوط هوایی و شرکت‌های اجاره ماشین معمولاً در کنار بزرگراه‌های متنه‌ی به فرودگاهها دیده می‌شوند، جایی که احتمالاً در بین رانندگان این راهها، استفاده‌کنندگان این خدمات وجود خواهند داشت.

برای به کاربردن داده‌کاوی در این مسئله، باید اول مشتری احتمالی خوب را تعریف کرده و سپس قواعد را یافت که اجازه هدف‌گیری مشتریان دارای خصوصیات مورد نظر را می‌دهد. برای اکثر شرکتها، اولین قدم به منظور استفاده از داده‌کاوی برای تشخیص مشتریان احتمالی خوب، ساختن یک مدل پاسخ<sup>۳</sup> است.

### انتخاب کانال ارتباطی

یافتن مشتریان احتمالی نیاز به ارتباط دارد. شرکتها به عمد از راههای مختلفی با مشتریان احتمالی تماس می‌گیرند. یکی از راه‌ها، روابط عمومی است که هدف آن تشویق رسانه‌ها برای پوشش داستانی درباره شرکت و نیز انتشار پیامهای مثبت افواهی می‌باشد. این کار برای برخی شرکتها بسیار مؤثر است ولی پیامهای بازاریابی مستقیم با روابط عمومی متفاوتند.

در اینجا تبلیغات و بازاریابی مستقیم مورد نظر است. تبلیغات می‌توانند روی جلد مجله، پنجره‌های باز شده مزاحم در برخی سایتهای تجاری، زیرنویسهای تلویزیونی در حین وقایع ورزشی مهم یا تبلیغات محصول در فیلمها باشد. این نوع تبلیغات، گروه‌های مردم را بر اساس

<sup>۱</sup>- Loyal

<sup>۲</sup>- Targeting

<sup>۳</sup>- Response Model

صفات مشترک هدف‌گیری می‌کند ولی پیام را برای هر فرد مشخص نمی‌کند. در بخش‌های آتی، از طریق تطابق پروفایل<sup>۱</sup> مشتریان احتمالی با پروفایل یک ناحیه جغرافیایی، به‌منظور انتخاب محل مناسب تبلیغ بحث می‌شود.

### پروفایل مشتری

اگر مشخصات عمومی مشتری مانند سن، جنسیت و آدرس با مشخصات مشتریان دیگر مقایسه شود، شرکت را قادر به شناسایی نوع افرادی می‌کند که محصولاتش را می‌خرند. این مشخصات به شرکت کمک می‌کند که محصولات دیگری برای همان گروه تولید کند یا استراتژیهای متفاوتی برای فروختن همان محصولات به بازارهای هدف دیگر توسعه دهد.

بازاریابی مستقیم، اجازه می‌دهد پیام برای هر فرد، شخصی شود. این کار می‌تواند از راه تماس تلفنی (مثل SMS)، پست الکترونیکی، کارت پستان و یا کاتالوگ رنگی گلاسه باشد. داده‌کاوی می‌تواند به تعیین کانالهای مؤثر برای هر گروه از مشتریان احتمالی کمک کند.

### انتخاب پیامهای مناسب

حتی برای فروختن محصول یا خدمتی یکسان، پیامهایی متفاوت برای افراد مختلف مناسب است. مثلاً، ممکن است یک روزنامه برای یکی به دلیل پوشش اخبار ورزشی و برای دیگری به خاطر پوشش اخبار سیاسی یا هنری جذاب باشد. وقتی محصول دارای تنوع بوده یا چند محصول پیشنهاد می‌شود، انتخاب پیام مناسب مهم‌تر هم می‌شود.

حتی در یک محصول نیز پیام مهم است. یک مثال کلاسیک، تعامل میان معیار هزینه و معیار راحتی است. برخی مردم به قیمت خیلی حساس بوده و تمایل دارند از تعاونی خرید کرده، نصفه شب تلفن کنند، بین مسیر هوایپما عوض کنند و سفرهایشان شامل شب آخر هفته باشد. دیگران حاضرند برای خدمات راحت‌تر، مبالغ اضافه بپردازنند. پیامی بر مبنای قیمت کمتر، نه

فقط در انگیزش راحت‌طلبان شکست می‌خورد بلکه ممکن است خطر راندن آنان به سمت محصولات کم‌سودتر را داشته باشد، درحالی‌که آنان تمایل دارند پول بیشتری خرج کنند. مدل‌های پاسخ که شامل یک فعالیت تبلیغی<sup>۱</sup> هستند، می‌توانند با هم ترکیب شده تا بهترین پیشنهاد را به مشتری بدهند. برای دسته‌بندی مشتریان به بخش‌های دارای تشابه فکری در پاسخگویی به پیشنهادات، می‌توان از فیلتر کردن مشارکتی استفاده کرد.

#### ۱۰-۲-۱- داده‌کاوی برای انتخاب محل مناسب تبلیغ

یک راه هدف‌گیری مشتریان احتمالی، نگاه کردن به مشتریان فعلی است. برای مثال یک نشریه کشوری از طریق پرسشنامه، مشخصات زیر را برای خوانندگانش به دست آورده است:

- ۵۸٪ تحصیلات عالی داشتند.

- ۴۶٪ مشاغل تخصصی یا مدیریتی داشتند.

- ۲۱٪ درآمد خانواری بیش از ۸ میلیون تومان در سال داشتند.

- ۷٪ درآمد خانواری بیش از ۱۲ میلیون تومان در سال داشتند.

درک این پروفایل از دو طریق می‌تواند به نشریه کمک کند. اول اینکه با هدف‌گیری مشتریان احتمالی مطابق با این پروفایل، می‌توان نرخ پاسخ به فعالیتهای ترویجی<sup>۲</sup> خود را افزایش داد. دوم، می‌توان فضای تبلیغاتی نشریه را به شرکهای علاقه‌مند به این نوع خوانندگان تحصیل کرده و پردرآمد فروخت. از آنجا که موضوع این بخش، هدف‌گیری مشتریان احتمالی است، باید ببینیم چگونه این نشریه از این پروفایل برای مرکز کردن فعالیتهای مشتری‌یابی خود استفاده کرده است. ایده اصلی ساده است. وقتی نشریه می‌خواهد در روزنامه آگهی بدهد، باید به دنبال روزنامه‌هایی باشد که شنوندگانشان مطابق این پروفایل باشند. باید در جاهایی کارتهای معرفی خود را روی پیشخوان مغازه‌ها بگذارد که منطبق بر این پروفایل باشند. وقتی می‌خواهد بازاریابی با پیامک (SMS) انجام دهد باید با مردمی مطابق این پروفایل تماس بگیرد. چالش داده‌کاوی، ارائه تعریف مناسبی از معنای انطباق با این پروفایل است.

<sup>۱</sup>- Campaign

<sup>۲</sup>- Promotional

### چه کسی با این پروفایل مطابقت دارد؟

یک راه تعیین تطابق مشتری با پروفایل مورد نظر، اندازه‌گیری شباهت یا فاصله آن دو است. بسیاری از فنون داده‌کاوی از ایده اندازه‌گیری شباهت به عنوان فاصله استفاده می‌کنند. برای مثال، مدل استدلال بر مبنای حافظه<sup>۱</sup>، روشی برای دسته‌بندی مشاهدات بر پایه دسته‌های مشاهده شده‌ای می‌باشد که در همان همسایگی وجود دارند. کشف خوش‌خودکار، فن داده‌کاوی دیگری است که بر اساس امکان محاسبه فاصله بین دو مشاهده برای یافتن گروه مشاهدات شبیه بهم کار می‌کند. در این مثال، به دنبال تعریف معیار فاصله‌ای برای درجه تطابق مشتری احتمالی با پروفایل موجود هستیم، داده‌ها شامل نتایج پرسشنامه بوده و نشان‌دهنده مشترکین در یک مقطع زمانی می‌باشند. در اینجا با این پرسشها روی رو هستیم: چه نوع معیارهایی مناسب این داده‌ها هستند؟ پروفایل‌ها چه نیازهایی را برآورده می‌کنند؟

دو فرد شرکت‌کننده در تحقیق را در نظر بگیرید. مریم، تحصیل کرده بوده، ۹ میلیون تومان در سال درآمد داشته و متخصص است. داود دپلمه بوده و ۴ میلیون تومان در سال درآمد دارد. کدام یک بیشتر با پروفایل خوانندگان تطابق دارند؟ جواب، بسته به نحوه مقایسه متفاوت است. جدول (۱۰-۱) راه مناسبی را برای محاسبه اختیار از روی پروفایل و معیار فاصله نشان می‌دهد.

جدول (۱۰-۱) محاسبه اختیارهای تطابق برای افراد با مقایسه معیارهای دموگرافیک

خواستگان نشانیه	تحصیل کرده	متخصص یا مدیر	درآمد بیشتر از ۸	درآمد بیشتر از ۱۲	جمع
امتیاز داود	امتیاز مریم	داود	مریم	خیر	بله
۰,۴۲	۰,۵۸	خیر	بله	۰,۴۲	۰,۵۸
۰,۵۴	۰,۴۶	خیر	بله	۰,۵۴	۰,۴۶
۰,۷۹	۰,۲۱	خیر	بله	۰,۷۹	۰,۲۱
۰,۹۳	۰,۹۳	خیر	بله	۰,۹۳	۰,۰۷
۲,۶۸	۲,۱۸				

این جدول امتیازی بر پایه نسبت توافق مخاطب با هر خصوصیت را محاسبه می‌کند. برای مثال چون ۵۸٪ خوانندگان تحصیل کرده‌اند، مریم امتیاز ۰,۵۸ برای این خصوصیت می‌گیرد. داوود که دیپلمه است امتیاز ۰,۴۲ می‌گیرد، زیرا ۴۲٪ خوانندگان نیز تحصیلات عالی ندارند. این امتیاز برای هر خصوصیت محاسبه شده و امتیازها جمع می‌شوند. امتیاز مریم ۰,۱۸ و امتیاز داوود ۰,۶۸ می‌شود. امتیاز بالاتر داوود نشان می‌دهد او از مریم به خوانندگان فعلی شبیه‌تر است. مشکل این روش این است که با اینکه طبق این امتیازدهی داوود مناسب‌تر از مریم به نظر می‌رسد ولی در واقع مریم به مخاطبین هدف نشریه یعنی افراد تحصیل کرده و پردرآمد شبیه‌تر است.

موفقیت این هدف‌گیری از مقایسه پروفایل خوانندگان با خصوصیات جمعیت‌شناسی<sup>۱</sup> معلوم می‌شود. این موضوع لزوم برخوردي پخته‌تر با اندازه‌گیری تطابق فرد با مخاطبین نشریه از طریق در نظر گرفتن خصوصیات جمعیت عمومی علاوه بر خصوصیات خوانندگان را ایجاب می‌کند. این روش، درجه تفاوت یک مشتری احتمالی از جمعیت عمومی و تفاوت خواننده از جمعیت عمومی را محاسبه کرده و سپس شباهت این دو تفاوت را اندازه می‌گیرد.

### دموگرافی

مطالعه آماری جمعیت (جمعیت شناسی) شامل خصوصیاتی مانند توزیع جغرافیایی، محیط فیزیکی، بیماری، ترکیب جنسیتی و سنی، و نرخ تولد و مرگ.

خواننده نشریه در مقایسه با جمعیت عمومی تحصیل کرده‌تر، متخصص‌تر و پردرآمدتر است. در جدول (۱-۱۰) ستونهای شاخص از تقسیم درصد خوانندگانی که ویژگی خاصی دارند بر درصد جمعیتی که دارای این ویژگی است به دست آمده‌اند. می‌توان خصوصیات خواننده را با جمعیت عمومی از طریق این شاخصها مقایسه کرد. دیده می‌شود که خواننده نشریه تقریباً سه برابر جمعیت عمومی تحصیل کرده است. به طور مشابه، خواننده تقریباً نصف جمعیت عمومی تحصیل نکرده است. با استفاده از شاخصها به عنوان امتیاز هر خصوصیت، مریم امتیاز

<sup>۱</sup>- Demographic

می‌گیرد در حالی که امتیاز داود فقط  $(۰/۹۵+۰/۲۱+۰/۴۰+۰/۸۶+۰/۲)/۴۲ = ۰/۴۲$  می‌شود. امتیازهای بر پایه شاخصها با توجه به جمعیت مخاطب هدف، متناسب‌تر می‌باشند. این امتیازها با معناتر هستند زیرا شامل اطلاعات اضافه تفاوت مخاطبین هدف با جمعیت عمومی می‌باشد.

### شاخصها به جای مقادیر خام

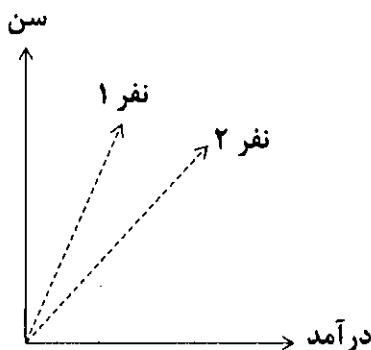
در مقایسه پروفایل مشتری، در نظر گرفتن پروفایل جمعیت عمومی مهم است. برای همین، استفاده از شاخصها اغلب از مقادیر خام بهتر است.

جدول ۲-۱۰) محاسبه امتیازها با در نظر گرفتن نسبتها در جمعیت

خبر				بله			
شاخص	خوانندگان جمهوری جمهوری جمهوری	خوانندگان جمهوری جمهوری جمهوری	شاخص	خوانندگان جمهوری جمهوری جمهوری	خوانندگان جمهوری جمهوری جمهوری	ناشر	
۰,۵۳	%۷۹,۹	%۴۲	۲,۸۶	%۲۰,۳	%۵۸	تحصیل کرده	
۰,۷۷	%۸۰,۸	%۵۴	۲,۴۰	%۱۹,۲	%۴۶	متخصص یا مدیر	
۰,۸۷	%۹۰,۵	%۷۹	۲,۲۱	%۹,۵	%۲۱	درآمدبیشتر از ۸	
۰,۹۵	%۹۷,۶	%۹۳	۲,۹۲	%۲,۴	%۷	درآمدبیشتر از ۱۲	

### مفهوم تشابه بر پایه زاویه و تفاوت بر پایه فاصله

می‌توان مفهوم تشابه را بر پایه تفاوت دو زاویه توضیح داد. هر ویژگی اندازه‌گیری شده یک بعد جدا در نظر گرفته می‌شود. با در نظر گرفتن مقدار متوسط هر ویژگی به عنوان مبدأ، پروفایل خوانندگان فعلی، برداری است که نشان می‌دهد خواننده نوعی، چقدر و در چه جهتی از جمعیت عمومی دور است. داده‌های یک مشتری احتمالی نیز یک بردار است. اگر زاویه بین این دو بردار کوچک باشد، آن گاه مشتری احتمالی نیز در همان جهت پروفایل خوانندگان با جمعیت عمومی تفاوت دارد.



شکل ۲-۱۰) شباهت (زاویه) و تفاوت (فاصله)

در شکل ۲-۱۰)، دو نفر (۱ و ۲) از نظر سن و درآمد مقایسه شده‌اند. می‌توان شباهت بین این دو را از طریق زاویه (کسینوس) بین دو بردار سنجید. ولی اگر یکی از آنها دقیقاً دو برابر دیگری سن و درآمد داشته باشد چه طور؟ در این صورت بردار یکی روی دیگری قرار می‌گیرد و شباهت ۱۰۰٪ می‌شود! پس معیار شباهت همیشه مناسب نیست و در بسیاری از موارد بهتر است از تفاوت دو نقطه ۱ و ۲ بر حسب معیار فاصله استفاده شود.

### نرمال کردن

به نظر شما آیا اندازه‌گرفتن درآمد بر حسب توانمندی یا ریال باید تأثیری در شباهت یا تفاوت دو نفر بگذارد؟ برای جلوگیری از این مشکلات معمولاً ابتدا داده‌های هر مشخصه مانند درآمد طوری نرمال می‌شوند که در یک مقیاس قرار گیرد. برای مثال درآمد کلیه افراد را تقسیم بر حداقل درآمد ممکن می‌کنیم تا کلیه درآمدهای بین صفر تا یک قرار گیرند.

### مشخصه‌های غیر عددی

چگونه مشخصه‌ای مانند جنسیت را در شباهت در نظر بگیریم؟ کافی است بنا بر قرارداد بگوییم اگر دو نفر هم‌جنس باشند، تفاوت آنها صفر و اگر دارای جنسیت مخالف باشند تفاوت آنها یک است. در این صورت با توجه به اینکه مشخصه‌های عددی را نیز نرمال کردایم (در فاصله صفر-یک) می‌توان همه مشخصه‌ها را در یک مقیاس در تفاوت یا شباهت در نظر گرفت.

### اندازه‌گیری مطابقت برای گروه‌های مشتریان

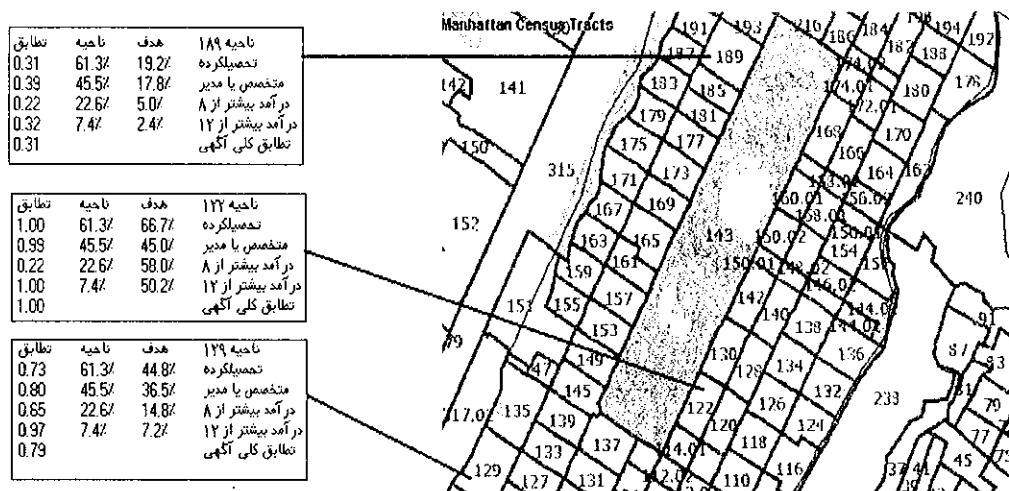
محاسبه امتیازهای بر پایه شاخص، قابل توسعه به گروه‌های بزرگ‌تر مردم است. اهمیت این موضوع در آن می‌باشد که ممکن است ویژگی خاص مورد استفاده برای هر مشتری یا مشتری احتمالی در دسترس نباشد. خوشبختانه و نه بر حسب تصادف، خصوصیات مطرح شده قبلی همگی خصوصیات جمعیت‌شناسی موجود در سرشماری آماری بوده و قابل اندازه‌گیری بر حسب نواحی جغرافیایی مانند منطقه سرشماری می‌باشند.

#### داده‌های هر منطقه سرشماری

یکی از فروض بازاریابی بر پایه ضرب المثل قدیمی "کبوتر با کبوتر، باز با باز، کند هم جنس با هم جنس پرواز." می‌باشد. یعنی مردم دارای علایق و سلیقه‌های مشابه در نواحی مشابه زندگی می‌کنند (داوطلبانه یا به دلیل الگوهای تاریخی تمايز). با توجه به این فرض، معامله با مردمی که قبلاً در میانشان و یا در نواحی مشابه مشتری داشته‌اید، ایده خوبی است. آمار سرشماری، هم برای یافتن مکانهای تمرکز مشتری و هم برای تعیین پروفایل نواحی مشابه، بالارزش است.

فرایند مورد نظر، نرخ‌گذاری هر منطقه سرشماری با توجه به تطابق آن با نشریه می‌باشد. ایده اصلی، تخمین نسبت هر منطقه سرشماری مطابق با پروفایل خوانندگان نشریه است. برای مثال اگر یک منطقه سرشماری دارای جمعیت بزرگسالی با ۵۸٪ تحصیل‌کرده باشد باز هم امتیاز مریم، ۱ یعنی تطابق کامل است. اگر فقط ۵۰٪ تحصیل‌کرده باشند، آنگاه امتیاز تطابق این خصوصیت ۰,۱ است. امتیاز تطابق کل، متوسط امتیازهای هر خصوصیت است.

شکل (۳-۱۰) مثالی از سه منطقه سرشماری را نشان می‌دهد. هر منطقه دارای نسبت متفاوتی از چهار خصوصیت مورد نظر است. این داده‌ها را می‌توان برای محاسبه امتیاز تطابق کلی هر منطقه ترکیب نمود. توجه کنید که همه ساکنان آن منطقه، امتیاز یکسانی می‌گیرند. این امتیاز بیانگر نسبتی از جمعیت آن منطقه است که مطابق پروفایل مورد نظر می‌باشد.



شکل ۳-۱۰) مثالی از محاسبه تطابق با خوانندگان در سه منطقه سرشماری مانهاتن

### استفاده از مشتریان فعلی برای یادگیری در مورد مشتریان احتمالی

یک راه مناسب برای یافتن مشتریان احتمالی، نگاه به همان مکانهایی است که بهترین مشتریان فعلی از آن جا می‌آیند. بنابراین باید راهی برای تشخیص بهترین مشتریان فعلی داشت. لازمه این کار، نگهداری سابقه نحوه دستیابی به مشتریان فعلی و وضعیت آنان در زمان دستیابی می‌باشد.

البته تکیه به مشتریان فعلی برای یادگیری در مورد مشتریان احتمالی، خطر انعکاس تصمیمات بازاریابی گذشته را دارد. مطالعه مشتریان فعلی، پیشنهادی برای جستجوی مشتریان احتمالی در مکانهای جدید نمی‌دهد. با این حال، عملکرد فعلی راه مناسبی برای ارزیابی کانالهای مشتری‌یابی موجود است. برای یافتن مشتریان احتمالی، مهم است بدانیم وقتی مشتریان فعلی، خودشان زمانی مشتری احتمالی بوده‌اند، چگونه به نظر می‌رسیدند. به طور ایده‌آل باید:

- ردگیری مشتریان را قبل از مشتری شدن شروع کرد.
- اطلاعات مشتریان جدید را در حین مشتری شدن آنان جمع کرد.
- رابطه بین داده‌های زمان مشتری شدن و نتایج مورد نظر آتی را مدل کرد.

- بخش‌های بعدی، این بحث را تکمیل می‌کنند.

### شروع ردگیری مشتریان قبل از مشتری شدن

خوب است اطلاعات مشتریان احتمالی قبل از مشتری شدن، ثبت شود. سایتهاي وب می‌توانند با ایجاد یک کوکی<sup>۱</sup> در اولین بازدید مشتری از سایت و سپس پروفایل کردن این مشتری بی‌نام با ثبت کارهایی که انجام می‌دهد، به این منظور برسند. وقتی بازدیدکننده از طریق همان برنامه پوشگر وب مانند *IE* و روی همان کامپیوتر قبلی، دوباره به سایت رجوع می‌کند، این کوکی شناخته شده و پروفایل به روز می‌شود. وقتی این بازدیدکننده تبدیل به مشتری یا کاربر ثبت‌نام کرده می‌شود، فعالیت منجر به انتقال بخشی از سابقه مشتری می‌گردد.

ردگیری پاسخها و پاسخ‌دهندگان در دنیای غیر وی‌بی نیز کار خوبی است. اولین اطلاعات مهم برای ثبت، پاسخ یا عدم پاسخ مشتری می‌باشد. مشخصات کسانی که پاسخ داده‌اند و آنهایی که پاسخ نداده‌اند، جزو لاینک مدل‌های پاسخ است. در صورت امکان، داده‌های اقدام بازاریابی محرك پاسخ، کanal دریافت پاسخ و زمان پاسخ را نیز باید ثبت کرد.

تشخیص پیام بازاریابی محرك پاسخ از میان پیامهای متعدد، می‌تواند دشوار و یا غیرممکن باشد. برای آسان کردن این تشخیص، فرمهای پاسخ و کاتالوگها دارای کد‌های تشخیص مناسب هستند. وب سایتها آدرس سایتی را که از آن مراجعه شده، ثبت می‌کنند. حتی می‌توان عملیات تبلیغات را از طریق تلفنهای جدا، بسته‌های پستی و یا آدرس سایتهاي مختلف تفکیک نمود. بسته به طبیعت محصول یا خدمت، ممکن است پاسخ‌دهندگان ملزم به ارائه اطلاعات اضافه در فرم محصول یا ثبت نام باشند. اگر انجام خدمت مستلزم داشتن اعتبار باشد، اطلاعات اعتبار درخواست می‌شود. اطلاعات جمع‌آوری شده در ابتدای ارتباط با مشتری از هیچ گرفته تا معاینه کامل پژوهشکی برای بیمه عمر متغیر است. اکثر شرکتها جایی در وسط هستند.

### جمع‌آوری اطلاعات از مشتریان جدید

وقتی یک مشتری احتمالی برای اولین بار مشتری می‌شود، فرصتی طلایی برای جمع‌آوری اطلاعات بیشتری پیش می‌آید. داده‌های مربوط به مشتری احتمالی قبل از تبدیل شدن به مشتری،

از نوع جغرافیایی و یا دموگرافیک هستند. بعید است لیستهای خریداری شده شامل چیزی غیر از نام، آدرس تماس و منبع لیست باشند. با داشتن آدرس می‌توان اطلاعات دیگری در مورد مشتریان احتمالی بر پایه خصوصیات همسایگان به دست آورد. با داشتن نام و آدرس با هم می‌توان اطلاعات سطح خانوار را از فراهم‌کنندگان داده‌های بازاریابی خرید. این نوع از داده‌ها برای هدف‌گیری بخش‌های عمومی مانند «مادران جوان» یا «نوجوانان حومه» مفید است ولی برای ایجاد رابطه فردی با مشتری به اندازه کافی دارای جزئیات نیست.

مفیدترین مشخصات قابل جمع‌آوری برای داده‌کاوی آتی، تاریخ خرید اولیه، کanal خرید اولیه، پیشنهاد پاسخ، محصول اولیه، امتیاز اعتبار اولیه، مدت تا پاسخ و محل جغرافیایی است. این مشخصات در عمل پیش‌بینی کننده خوبی برای مواردی مانند مدت مورد انتظار رابطه، بدحسابی<sup>۱</sup> و خریدهای اضافی هستند.

#### پیش‌بینی نتایج آتی بر اساس متغیر زمان مشتری شدن

با ثبت همه اطلاعات ممکن در زمان مشتری شدن و سپس ردگیری مشتریان در طول زمان، می‌توان از داده‌کاوی برای مرتبط کردن متغیرهای زمان خرید به نتایج آتی مانند طول عمر مشتری، ارزش مشتری و ریسک پیش‌فرض استفاده کرد. این اطلاعات را می‌توان برای هدایت تلاشهای بازاریابی و تمرکز روی کانال‌ها و پیامهای مؤثر (دارای بهترین نتایج) به کار برد. کشف اینکه برخی کانال‌ها، مشتریانی با طول عمر دو برابر کانال‌های دیگر می‌باشد، غیر معمول نیست. با فرض امکان تخمین ارزش ماهانه مشتری و داشتن طول عمر، می‌توان ارزش ریالی مشتری هر کanal را محاسبه کرد. برای نرخ گذاری کانال‌ها، ارزش ریالی مشتری به اندازه هزینه تماس، مهم است.

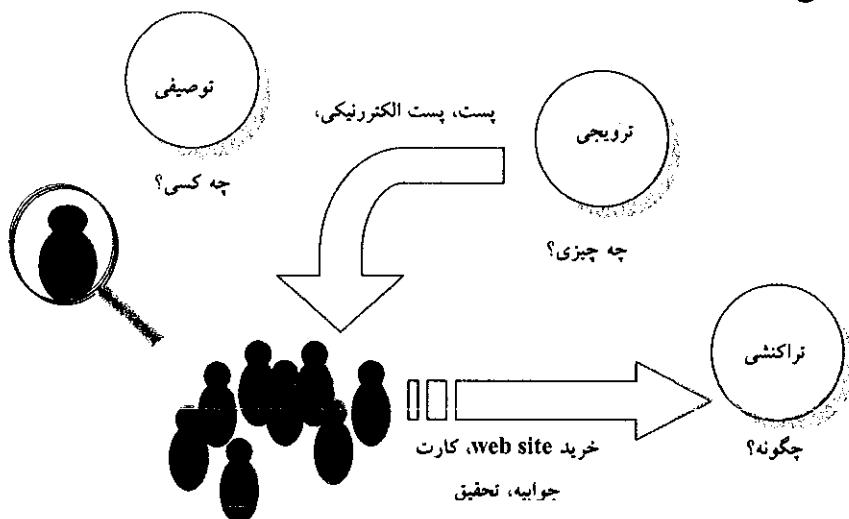
#### ۱۰-۱-۳- داده‌های مشتریان

اگر بخواهید از مشتری تصویر درستی داشته باشد باید از جهات مختلف به مشتری بنگرید و همه اطلاعات درباره مشتری را از بخش‌های مختلف سازمان در کنار هم بگذارید، در این صورت است که این تصویر می‌تواند خصوصیات و رفتار صحیح مشتری را نمایان سازد. برای

مثال چون بخش فروش با مشتری ارتباط برقرار می‌کند، تصویری از نیازهای مشتری را دریافت می‌کند. زمانی که مشتری برای حل مشکلی با بخش خدمات تماس می‌گیرد در واقع تصویر دیگری از خود را به سازمان نمایان می‌سازد (رضایت یا عدم رضایت). هر کدام از این نقاط تماس<sup>۱</sup> فرصتی برای تعامل هستند و سازمان می‌تواند از طریق آنها با مشتریانش در تماس باشد یا مشتریان با سازمان در تماس باشند، به همین دلیل به آنها نقاط تماس گفته می‌شود. در واقع نقاط تماس بهترین اطلاعات را درباره مشتریان تامین می‌کنند. ولی چون این اطلاعات به صورت پراکنده و مجزا وجود دارد، هر قسمت به تنها برای ارائه یک تصویر کامل از مشتری کافی نیست. درنتیجه یکپارچگی این داده‌ها که از منابع مختلفی می‌آیند در ارائه تصویر کامل از مشتری بسیار حائز اهمیت است [۱].

داده‌های مرتبط با مشتریان را می‌توان از منابع مختلفی به دست آورد شکل (۱۰-۴)، در سیستم‌های داده‌کاوی مدیریت ارتباط با مشتری سه نوع اصلی داده وجود دارد که عبارتند از:

- توضیح اینکه مشتری کیست.
- توضیح اینکه چه بازاریابی یا تبلیغ فروشی برای مشتری انجام شده است.
- توضیح اینکه مشتری در مقابل این تبلیغات چه واکنشی نشان داده است.



شکل ۱۰-۴) برای فراهم کردن اطلاعات کافی برای داده‌کاوی باید به دنبال چه کسی، چه چیزی و چگونه باشید.

اگر شما این سه چیز درباره مشتریان یا حتی درباره افرادی که هنوز مشتری شما نیستند بدانید، در واقع داده کافی برای شروع پیش‌بینی را دارید. شما می‌توانید با کمک داده‌کاوی الگوهایی را از میان داده‌ها کشف کنید یا حتی با بهره‌گیری از این تجربیات، تراکنش‌های بازاریابی و فروش با این مشتریان را بهینه کنید. بدون دانستن اینکه مشتری کیست، چه کاری انجام می‌دهد و واکنش آن چیست، بهینه‌سازی یا بهبود سیستم امکان‌پذیر نمی‌باشد.

برای داشتن سودآورترین تراکنش با مشتری تا جای ممکن و بهینه کردن عملکرد سیستم CRM، باید توانایی تفکیک مشتریان خوب و بد و مشتریان سودآور و غیرسودآور را داشته باشید. باید بدانید آنها چه کسانی هستند و چه تفاوت‌هایی با هم دارند. به‌طور مشابه برای اینکه بدانید در تبلیغات و بازاریابی چگونه سرمایه‌گذاری کنید، نیاز به دانستن کارهایی که برای هر مشتری انجام می‌دهید و نتایج پیگیری این کارها دارید. شما باید تعداد زیادی تجربه‌های کوچک در امور تبلیغاتی و بازاریابی با مشتریان اصلی تان داشته باشید، و توجه کنید که تفاوت در هر تجربه بهترین راه برای پی بردن به کارهایی که باید انجام دهید و نباید انجام دهید، می‌باشد. برای یافتن ارزش واقعی سیستم، باید بتوانید نتایج را اندازه‌گیری کنید. اگر مشخص نشود که نتیجه تجربه خوب یا بد بوده است، پس واقعاً چیز جدیدی که بتواند برای بهبود سیستم در دفعات بعد استفاده شود، حاصل نشده است. گروه‌بندی دسته‌های مختلف تجربه‌ها نیز مفید است چون باعث شکسته شدن داده‌های ذخیره شده در بانک اطلاعاتی به انواع مختلف داده می‌شود. همچنین این دسته‌بندیها توصیف خوبی از منابع داده ارائه می‌دهند. قراردادن داده در این سه دسته به شما برای داده‌کاوی موفق و تولید یک سیستم بهینه CRM کمک شایانی خواهد کرد.

### داده توصیفی

داده توصیفی شامل توضیحاتی درباره مشتری یا مصرف‌کننده است. این نوع داده‌ها معمولاً به‌طور خلاصه در ستونهای مختلف جدول اطلاعات مشتریان در بانک اطلاعاتی ذخیره می‌شود. این‌گونه داده‌های توصیف‌کننده مشتری شامل اطلاعاتی مثل سن، جنسیت، موقعیت منزل، تعداد فرزندان، درآمد خانگی و درآمد فردی هستند. این اطلاعات قابل تغییر هستند ولی معمولاً زودتر از یکسال عوض نمی‌شود. البته اطلاعاتی مانند آدرس و تلفن نیاز به بروز رسانی فصلی یا

حداقل شش ماه یکبار در بانک اطلاعاتی دارند. این داده‌ها شامل موارد کلی ذیل هستند: دموگرافی، مالی و پروفایل.

### رفتار مشتری است که مهم می‌باشد

هدف داده‌کاوی در CRM، اصلاح نسبت رفتارهای مثبت به رفتارهای منفی است. اطلاعات ایستا یا همان توصیفی مانند سن یا کُپستی صرفاً جانشینی برای اطلاعات واقعی مهم یعنی رفتار مشتری هستند. پس چرا فقط از داده‌های رفتاری مانند سابقه خرید مشتری استفاده نمی‌کیم؟ زیرا در بسیاری از اوقات اطلاعات کافی از سابقه رفتاری مشتریان موجود نیست در حالی که یک مشخصه ایستا مانند کد پستی مشتری، می‌تواند پیش‌بینی کننده خوبی برای بسیاری از رفتارهای مشتری باشد.

### داده تبلیغاتی

داده تبلیغاتی شامل اطلاعات کارها و فعالیتهایی است که برای مشتری صورت گرفته است. قدرت این نوع داده معمولاً به پیچیدگی سیستم CRM بستگی دارد. این داده‌ها شامل موارد زیر است:

- لیستی از کارهایی است که برای تبلیغات صورت گرفته است مثل پست، کاتالوگ، نمونه‌های تبلیغاتی یا کارتهای تخفیف
- تبلیغات تعاملی مثل تبلیغ در تلویزیون، رادیو، روزنامه، و مجله‌های تبلیغاتی و غیره
- اطلاعات دقیقی مثل ارسال پست الکترونیکی و تعداد کلیکهای کاربران قابل شناسایی در سایتها و وب

انواع اطلاعاتی که می‌تواند جمع‌آوری شود عبارت است از:

- نوع تعامل<sup>۱</sup>: فروش، بازاریابی از راه دور، تبلیغات چاپی، تبلیغات رسانه‌ای، تبلیغات وبی.
- توصیف تعامل: مانند رنگ کارت پستال.

- رسانه: بازاری که تبلیغات در آن صورت می‌گیرد، وب سایتهايي که آگهی تبلیغاتی در آنها قرار دارد.
- زمانبندی: زمان تعامل.
- توصیف قصد و نیت: یک توصیف کامل از اینکه برای چه کسی تعامل معنی دارد و چرا؟ (مثلًاً چرا این رنگ یا موسیقی زمینه باید انتخاب شود.)
- مالی: هزینه‌های ثابت و متغیر تعامل.

### داده تراکنشی

به طور کلی داده‌های تراکنشی، داده‌هایی است که مربوط به تعامل با مشتریان می‌باشد. این داده‌ها می‌توانند هر چیزی از یک تماس تلفنی برای درخواست خدمات گرفته تا توضیح و توصیف محصولاتی که مشتری خریده است باشد. این داده‌ها هم مثل داده‌های تبلیغاتی، می‌توانند خیلی سریع در طول زمان تغییر کند. بنابراین طبیعی است که باید داده‌ها در ساختاری ذخیره شود که به سادگی قابل به روز رسانی و تغییر باشد. اطلاعات تراکنشی با اطلاعات توصیفی مشتریان که در طول زمان اساساً تغییر نمی‌کند متفاوت است.

تغییر آرایش داده‌های تراکنشی در فاصله زمانی کوتاهی می‌تواند خیلی چشمگیر باشد. برای مثال معرفی محصولات جدید یا مورد توجه قرار گرفتن محصولات قدیمی و فروش بیشتر آنها، الگوی محصولات فروخته شده را تغییر می‌دهد. این داده‌ها شامل موارد ذیل است: خرید، کلیک صفحه وب، تماس تلفنی، پست الکترونیک، بازدید از مغازه و پست فیزیکی.

### لزوم تجمیع داده‌های تراکنشی

بسیاری از روش‌های داده‌کاوی نیاز به یک رکورد اطلاعاتی برای هر نمونه آموزشی (در CRM، هر مشتری) دارند. داده‌های ایستایی مانند سن و جنسیت برای هر مشتری فقط یک عدد در هر رکورد هستند. در حالی که داده‌های واقعه‌ای یا همان تراکنشی مانند سابقه خرید برای هر مشتری ممکن یک یا چند رکورد می‌باشند. برای همین لازم است این داده‌ها از نظر زمانی یا مقداری تجمیع شوند. برای مثال محاسبه شود که هر مشتری از زمان شروع خریدش تا کنون، هر سه ماه چه مبلغی خرید داشته است یا از هر محصول چه تعداد خریداری کرده است. هر مشخصه تراکنشی غیر عددی (مانند نوع محصول) یا ترتیبی، نامزد خوبی برای تجمیع است.

بسته به کاربرد، راههای مختلفی برای تعریف تجمعی وجود دارد (مثلاً ماهانه یا فصلی، متوسط خرید یا حداکثر مقدار خرید) بنابراین می‌توان از روی داده‌های واقعی‌ای تعداد زیادی مشخصه برای هر مشتری ساخت. این موضوع، مسئله بُعد زیاد داده و لزوم استفاده از فنون کاهش بعد را ایجاد می‌کند که قبلاً بحث شده است.

## ۱۰-۲- برخی کاربردهای داده‌کاوی در مدیریت ارتباط با مشتری

داده‌کاوی ابزاری بنیادی است که برای آشکارسازی خصوصیات جمعیت‌شناسختی مشتریان الزامی می‌باشد. از فنون داده‌کاوی می‌توان برای دستیابی به دامنه وسیعی از اهداف مختلف استفاده کرد. چند مثال از کاربردهای آن عبارتند از:

- شناسایی مشتریان سودآور و پروفایل آنان
- پیش‌بینی رفتار خرید مشتری
- رتبه‌بندی رویگردنی مشتری به منظور ارائه برنامه‌های مؤثر حفظ مشتری
- تمرکز تلاشهای بازاریابی بر مشتریان بالقوه‌ای که احتمال خرید کردن بیشتری دارند
- تخمین کارآمدی تبلیغات
- تخمین و اولویت‌بندی ریسک اعتباری مشتری
- برآورد میزان جدی بودن احتمالی مشتری
- فروش کناری<sup>۱</sup> و بالاسری<sup>۲</sup> به مشتریان بر اساس خرید محصولات قبلی
- هدفگیری مستقیم بازاریابی به سمت افرادی که بیشترین احتمال پاسخ را دارند
- پیش‌بینی کلاهبرداری و تقلبات
- بهینه‌سازی سهم سبد خرید مشتری

به طور خلاصه با به کارگیری انواع روش‌های داده‌کاوی، سازمان از مفهوم فروش صرف به سوی خدمت‌رسانی به مشتریان حرکت می‌کند. برخی از این کاربردها به اختصار توضیح داده می‌شوند.

<sup>۱</sup>- Cross-Selling

<sup>۲</sup>- Up-Selling

### مدلسازی حفظ و رویگردانی

حفظ بالارزش‌ترین مشتریان و دانستن اینکه کدام‌یک در خطر رویگردانی هستند می‌تواند به‌طور چشمگیری سودآوری سازمان را تحت تأثیر قرار دهد. سازمان در این مدل باید از موارد زیر آگاه باشد:

- اینکه کدام مشتریان در حال رویگردانی به سوی رقیب هستند و دلایل آن.

- اینکه کدام‌یک از ارزشمندترین مشتریان سازمان در خطر هستند.

- اینکه آیا سازمان بودجه حفظ مشتریان را برای با ارزش‌ترین مشتریان خرج می‌کند یا خیر.

- اینکه آیا سازمان راههای متناوب معنی‌دار و مؤثری به منظور تماس با مشتریان دارد یا خیر. از دست دادن مشتریان برای شرکتها تا حدی قابل اجتناب است. از طریق تجزیه و تحلیل داده‌ها می‌توان مدل‌هایی را به منظور پیش‌بینی احتمال خروج مشتریان و احتمال جذب آنان به دیگران از طریق تبلیغات فروش و عملیات تبلیغاتی ساخت.

با این مدل‌ها سازمان می‌تواند با تهیه مشخصه‌های افرادی که در گذشته رویگردان شده‌اند، مشتریانی را تعیین کند که دارای بیشترین گرایش به کاهش یا قطع ارتباطند. مدل حفظ مشتری، پتانسیل یک مشتری را در ماندن با سازمان پس از رخ دادن اتفاقات احتمالی، بررسی می‌کند. مدل رویگردانی، احتمال توقف خریدهای مشتریان فعال را بررسی می‌نماید. با این اطلاعات می‌توان سیاستهای پیشگیرانه‌ای در نظر گرفت و مشتریان را به منظور توجه ویژه به آنان، فعالانه تعقیب یا علامت‌گذاری نمود.

### مدلسازی پرهیز از ریسک

در این مدل سعی در پرهیز از کسب مشتریان غیرسودآور است. داده‌کاوی می‌تواند به‌طور تقریبی پیش‌بینی کند که کدام مشتریان بالقوه تبدیل به مشتریان بالفعل می‌شوند. ولی به‌طور خاص مفید است که تعیین شود کدام مشتریان سودآور خواهند بود. این موضوع مهم باید در نظر گرفته شود که به هر حال برخی مشتریان جدید بدھی خود را نخواهند پرداخت و شرکت مجبور به احتساب زیانهای وارده می‌باشد. این مدل، رفتار خرید، رفتار پرداخت، تاریخچه اعتباری و دیگر عوامل را بررسی می‌کند.

### مدلسازی فروش جانبی

فروش جانبی و فروش بالاسری در *CRM* اموری محوری محسوب می‌شوند. تعامل سازمان با مشتریان فرصتی اساسی برای بازاریابی محصولات یا خدمات اضافی ایجاد می‌کند. در این زمینه می‌توان داده‌های متفاوتی را در نظر گرفت مانند اینکه مشتریان چه چیزی خریداری می‌کنند، علایق آنها چیست، چه کالاهای یا خدماتی مد نظر شان است یا درباره کدام محصولات یا خدمات استعلام می‌کنند. در ابتدا باید در نظر گرفت که کدام محصولات و یا خدمات توان بالقوه‌ای در فروش جانبی دارند. این کار با قضاوت بر مبنای داده‌های فروش گذشته انجام می‌گیرد. به عبارت دیگر، تعیین می‌شود که یک مشتری منفرد کدام محصولات را بیشتر خریداری کرده است. سپس تعیین می‌گردد که پروفایل چه مشتریانی بیشترین تناسب را با گروه‌های متنوع محصول دارد. با این کار مشتریانی شناسایی می‌شوند که بیشترین تمایل را به خرید محصولات مشابه دارند ولی هنوز شرکت به این‌گونه مشتریان توجهی نشان نداده است. در این حالت سازمان فرصتی عالی برای تبلیغ آن محصولات به مشتریان فعلی خود پیدا می‌کند. تمام این فرایند در محدوده *CRM* می‌باشد.

### مدلسازی سودآوری

در این روش، ارزش طول عمر (*LTV*)<sup>۱</sup> مشتری اندازه‌گیری و تنظیم می‌شود تا معلوم گردد مشتریان در دوره طول عمر خریداری چه چیزی را خریداری می‌کنند یا پتانسیل خرید چه چیزهایی را دارند. این روش یکی از روش‌های تحلیلی داده‌کاوی است که از شیوه‌های محاسبه *LTV* و استفاده از آن برای دسته‌بندی کردن مشتریان استفاده می‌کند. فهم عناصر اصلی سودآوری کمک می‌کند تا سازمان درک نماید چه زمانی مشتریان سودآور خواهند شد و آیا اصلاً این اتفاق (سودآور شدن مشتریان) روی خواهد داد یا خیر.

### مدلسازی تجزیه و تحلیل اینترنتی

رفتار اینترنتی، یعنی چگونگی گشت و گذار مشتری درون صفحات اینترنتی مربوط به شرکت را می‌توان به منظور فهم رفتار و ترجیحات مشتری ضبط و تجزیه و تحلیل نمود.

ترجیحات مشتری شامل موارد زیر است: زمانی که مشتری صرف مشاهده یک صفحه می‌کند، کدام پیوندها را انتخاب می‌کند، به کدام آگهی‌ها توجه بیشتری دارد، از طریق کدام صفحه وارد سایت شرکت شده و از کدام صفحه خارج می‌گردد و نظایر آن.

همه این موارد به اطلاعات آماری پایگاه داده‌های شرکت تبدیل می‌شوند. این داده‌ها از این نظر مفیدند که به کمک آنها شرکتها می‌توانند مکانیزم فروششان را بشناسند و بدانند که چگونه محصولات و دیگر اقلام جانبی را که ممکن است مورد علاقه مشتریان باشد مکانیابی نمایند. این روش همچنین در صد ترک صفحه را نیز نشان می‌دهد. یعنی معلوم می‌کند که بازدیدکنندگان از روی کدام صفحات می‌پرند، از کدام صفحات اجتناب می‌کنند یا کدام صفحات موجب می‌شود که بازدیدکننده سایت وب شرکت را ترک کند.

تأکید این اطلاعات بر صفحاتی است که نیازمند ارزیابی بیشتر از نظر محظوظ و سادگی استفاده می‌باشند. البته با اینکه نرم‌افرازهای قوی تجزیه و تحلیل اینترنتی فراوانند ولی هیچ‌کدام از آنها قادر نیستند به تهایی تصویر کاملی از رفتارهای بازدیدکنندگان سایت ارائه دهند. مثلاً نمی‌توان به راحتی فهمید که آیا بازدیدکننده‌ها محصولات را برای خودشان می‌خرند یا به عنوان هدیه و برای دیگران می‌خرند، اینکه مشتریان سودآور هستند یا فقط در سایت به گشت و گذار می‌پردازند و غیره. بنابراین دموگرافی بازار هدف را نمی‌توان به تنهایی با استفاده از ابزار تجزیه و تحلیل اینترنتی انجام داد.

### **بازاریابی مستقیم**

هدف تبلیغات آن دسته از مشتریان احتمالی است که در مورد تک‌تک آنان اطلاعات کافی نداریم. بازاریابی مستقیم حداقل نیاز به مقداری اطلاعات اضافه مانند نام و آدرس یا تلفن یا پست الکترونیک دارد. در بسیاری از کشورها، داده‌های قابل توجهی درباره بخش بزرگی از جمعیت در دسترس است. قبل از برنامه‌ریزی برای استفاده در بازاریابی، لازم است دسترسی به داده‌ها در بازار مورد نظر و محدودیتهای قاعده‌ی استفاده، بررسی شود. مشکل این است که حتی از طریق غربالهای بدیهی، نسبت مشتریان مورد بررسی به مشتریان احتمالی پاسخ‌دهنده، بسیار زیاد باشند. بنابراین یکی از کاربردهای اصلی داده‌کاوی برای یافتن مشتریان احتمالی،

«هدف‌گیری» است. یعنی به دنبال یافتن آن مشتریان احتمالی هستیم که احتمال بیشتری دارد به یک پیشنهاد خاص پاسخ دهد.

فعالیتهای بازاریابی مستقیم عموماً نزخ پاسخی کمتر از ۱۰٪ دارند. این مدلها با تشخیص آن دسته از مشتریان احتمالی که احتمال پاسخ‌گیری به مشتری‌بایی مستقیم بیشتر است، نرخهای پاسخ را بهبود می‌بخشند. متفاوت‌ترین مدل‌های پاسخ، تخمینی واقعی از احتمال پاسخ می‌دهند. البته الزامی برای محاسبه احتمال واقعی پاسخ نیست بلکه داشتن مدلی که مشتریان احتمالی را به ترتیب بیشترین امتیاز پاسخ، رتبه‌بندی کند کافی است. با داشتن یک لیست مرتب شده می‌توان درصد پاسخ‌گویان در یک فعالیت بازاریابی مستقیم را با پست به افراد بالای لیست یا تعامل با آنها، زیاد کرد.

### ۱۰-۲-۱- لایه‌های کشف الگو

وقتی داده از اینباره داده استخراج شد و مرحله آماده‌سازی را طی کرد، هرکدام از روش‌های داده‌کاوی می‌تواند برای پاسخ به سوالات کسب وکاری که در ذهن است استفاده شود. دسته‌بندی، خوشبندی، رگرسیون، الگوهای مکرر<sup>۱</sup>، قواعد تلازمنی و سریهای زمانی مرسوم‌ترین روش‌های داده‌کاوی هستند. همان‌طور که در جدول (۱۰-۳) نشان داده شده کشف الگو شامل چهار لایه مهم است [۲].

جدول (۱۰-۳) لایه‌های کشف الگو

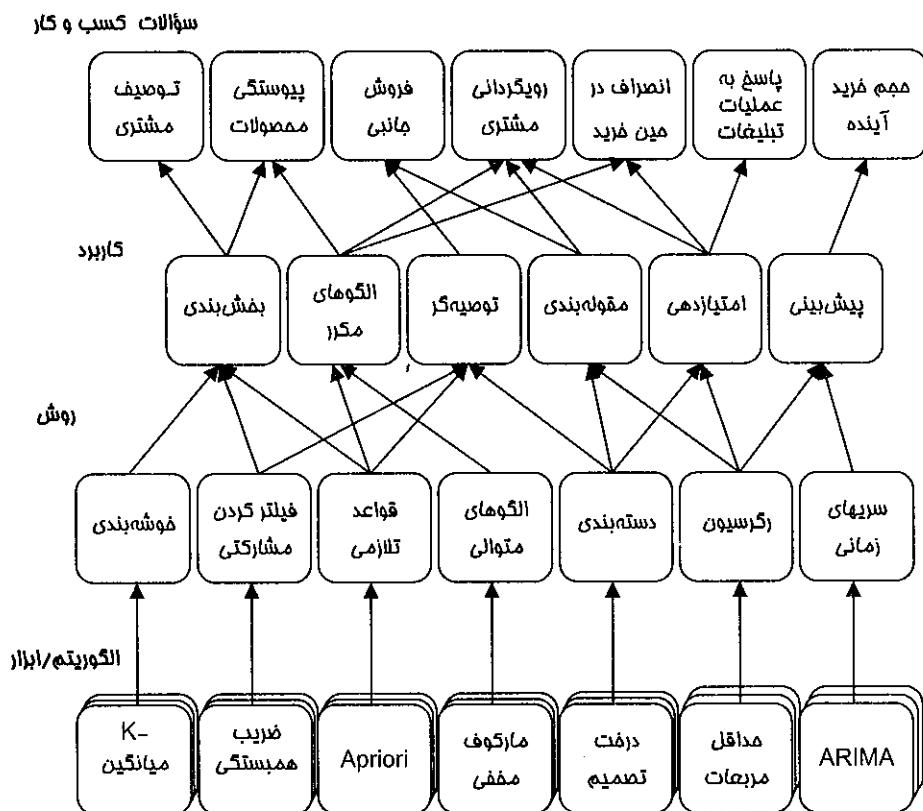
<p>سوالهای تجاری مانند توصیف مشتری کاربردها مانند امتیازدهی، پیش‌بینی برای حل یک نوع خاص سوال کسب وکار. روشها مانند سریهای زمانی، طبقه‌بندی با توجه به نوع داده خروجی که نتیجه فرآیند بر روی داده ورودی است <b>الگوریتمها</b></p>	<p>لایه اول لایه دوم لایه سوم لایه چهارم</p>
---	--

<sup>۱</sup>- Frequent Pattern

## سؤال استراتژیک و سؤال عملیاتی

سؤال استراتژیک، سؤالی است که جواب آن راهنمای تأثیرگذار بر یک تصمیم است. یک مثال می‌تواند تحلیل برای تعیین محرکهای کلیدی رضایت مشتری باشد. در مثال دیگر می‌توان با خوشبندی بر حسب کل دفعات خرید هر مشتری، مشتری با ارزش را تعریف کرد. سؤال عملیاتی، سؤالی است که نتایج آن مستقیماً برای افزایش وقایع سودآور استفاده می‌شود. برای مثال، در مدلسازی پاسخ بازاریابی مستقیم، خود امتیازهای مدل مهم‌ترین نتیجه هستند. اگر چه داشتن الگو یا مدل قابل توضیح، مطلوب است، ولی خود مدل هدف تحلیل می‌باشد.

در شکل (۵-۱۰) این لایه‌ها به طور جزئی تر نشان داده شده‌اند. البته امکان نمایش تمام حالات ممکن نمی‌باشد.



شکل (۵-۱۰) لایه‌های کشف دانش در مدیریت ارتباط با مشتری

### دسته‌بندی (در مدیریت ارتباط با مشتری)

آیا می‌توان از روی خصوصیات ظاهری، تشخیص داد که یک فرد زن است یا مرد؟ این سؤالی است که دسته‌بندی به دنبال جواب آن است. چگونه می‌توان این کار را کرد؟ ابتدا باید تعدادی زن و مرد داشته باشید که خصوصیات فیزیکی و الگوهای رفتاری آنها ثبت شده باشد. حالا اگر فرد جدیدی را ببینید، می‌توانید از روی شباهت خصوصیات او را به یکی از دو دسته زن یا مرد، تشخیص دهید. برخی اوقات به دلیل کامل نبودن اطلاعات و یا غیرعادی بودن فرد، امکان اشتباه نیز وجود دارد.

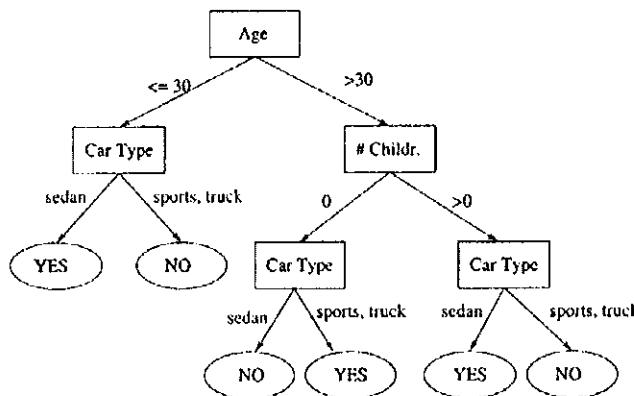
اگر دسته‌ها از قبیل مشخص باشند (مثلاً اینکه خرید کرده یا نکرده) می‌توان به کمک روش دسته‌بندی ابتدا به وسیله داده‌های موجود (داده‌های آموزشی) مدلی ساخته و از آن برای پیش‌بینی موارد جدید استفاده کرد.

### درخت تصمیم (در مدیریت ارتباط با مشتری)

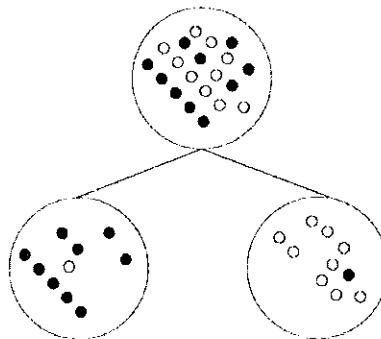
درخت تصمیم به دلیل سادگی بیان دانش یافته شده و قدرت پیش‌بینی مناسب، یکی از پرکاربردترین الگوریتمهای دسته‌بندی می‌باشد. دانش یافته شده به شکل قواعد اگر-آنگاه بیان می‌شود.

مسابقه ۲۰ سؤالی، نمونه‌ای از درخت تصمیم می‌باشد. در هر مرحله از مسابقه با پرسیدن یک سؤال، فضای جواب را محدودتر می‌کنیم تا در نهایت به جواب مشخصی که در ذهن طرف مقابل است برسیم. چه سؤالی را ابتدا انتخاب می‌کنید؟ سؤالی را ابتدا انتخاب می‌کنیم که فضای جواب را به خوبی به دو قسمت تقسیم کند.

مانند روش‌های دیگر دسته‌بندی، هدف در درخت تصمیم، تعیین دسته (مثلاً خرید یا عدم خرید مشتری) از روی مشخصات می‌باشد. برای اینکار ابتدا مشخصه‌ای (متغیری) انتخاب می‌شود که با گذاشتن شرط بیشتر یا کمتر از یک مقدار روی آن، بهتر از همه متغیرهای دیگر قدرت تفکیک کلیه رکوردهای موجود را داشته باشد.



برای مثال هدف در شکل بالا، پیش‌بینی اشتراک افراد در یک مجله است. با بررسی روی تک‌تک مشخصات مشتری معلوم شده است که بهترین مشخصه تفکیک کننده مشتری، سن مشتری و آستانه تصمیم آن سن ۳۰ سال است. طبق شکل ذیل با بررسی شرط بیشتر بودن سن از ۳۰ سال، ۲۰ فرد موجود را که ۱۰ نفر از آنها مشتری هستند، به دو دسته تقسیم می‌شوند که در یکی از ۱۰ نفر، ۹ نفر مشترک هستند و در دیگر از ۱۰ نفر، ۹ نفر مشترک نیستند.



این فرایند برای شاخه‌های زیرین درخت دنبال شده و هر بار مشخصه‌های دیگری انتخاب می‌شود که بهتر از بقیه، داده‌ها را تفکیک می‌کند. این روند تا رسیدن به حد قابل قبولی از خطا یا رسیدن به دسته‌های ۱۰۰٪ خالص ادامه می‌یابد.

یکی از موارد مهم در دسته‌بندی، استحکام می‌باشد. همان‌طور که متوجه شده‌اید امکان اشتباه در دسته‌بندی وجود دارد. بنابراین پس از ساختن مدل باید خطاهای دسته‌بندی را بررسی کرد تا در حد مقبولی باشند. ممکن است مدل برای داده‌های آموزشی خیلی خوب و کم خطأ دسته‌بندی

کند ولی در مورد داده‌های جدید (مثلاً تصمیم‌گیری در مورد مشتریان جدید) دارای خطای زیادی باشد. معیار ارتباط خطای آموزش به خطای استفاده در عمل، استحکام نام دارد. مدلی خوب است که دارای استحکام بالایی باشد.

### خوشه‌بندی (در مدیریت ارتباط با مشتری)

انسان به طور فطری تمایل به گروه‌بندی اشیاء و مفاهیم بر اساس شباهت یا تفاوت (فاسله) دارد. برای مثال کودک می‌آموزد که موجودات زنده دو دسته هستند: آنهایی که سبزند و راه نمی‌روند (گیاهان) و آنهایی که معمولاً قهوه‌ای‌اند، حرکت می‌کنند و احتمالاً خطرناکند (حیوانات)!

گروه‌بندی در ادبیات داده‌کاوی، خوشه‌بندی<sup>۱</sup> یا خوشه‌یابی نامیده می‌شود. روش خوشه‌بندی برای توصیف گروه‌های مختلف در مجموعه داده به کار می‌رود و بر خلاف دسته‌بندی خوشه‌ها از قبل مشخص نیستند. معمولاً ابتدا خوشه‌بندی انجام شده و بعد خوشه‌ها به عنوان نام دسته‌ها برای دسته‌بندی به کار می‌روند.

پس از خوشه‌بندی، برای تعبیر خوشه‌ها، نماینده هر خوشه را در نظر می‌گیرند. نماینده می‌تواند یکی از مشاهدات میانه خوشه و یا میانگین همه مشاهدات یک خوشه در نظر گرفته شود. با توجه به تفاوت و شباهت نماینده هر خوشه به نماینده کل داده‌ها از نظر مشخصات (مانند سن، جنسیت، ...) می‌توان هر خوشه را تعبیر کرد. مثلاً به یک خوشه نام «مادران جوان» و به خوشه دیگر نام «نوجوانان حومه» را داد. در CRM گروه‌بندی از آن جهت اهمیت دارد که ایجاد پروفایلی برای مشتریان مختلف را فراهم می‌کند و امکان برنامه‌ریزی استراتژیک روی گروه‌های مشتریان را می‌دهد. از طرف دیگر امکان انواع گروه‌بندی‌های دیگر مانند محصولات مشابه را ایجاد می‌کند.

ممکن است بسیاری از خوشه‌های یافته شده، بدیهی و فاقد دانش جدید و مفیدی برای ما باشند. بنابراین لازم است خوشه‌های یافته شده از نظر جالب بودن<sup>۲</sup> طبق نظر خبرگان بررسی شوند.

<sup>۱</sup>- Clustering

<sup>۲</sup>- Interestingness

گاهی اوقات، نشان دادن شباهت و تفاوت خوشة‌های مختلف با هم مفید می‌باشد. راه متداول برای اینکار شبکه‌های کوهونز یا همان *SOM* است.

### خوشبندی هدف‌گذاری شده (در مدیریت ارتباط با مشتری)

هر چند گروه‌های تشکیل شده بر اساس مشخصه‌هایی مانند سن، جنسیت یا رفتار قابل تمایز از هم هستند، ولی این ممکن است این گروه‌ها در بافت کسب و کار باعث نباشند. اگر همه گروه‌ها دارای عمر مشتری متوسطی بوده یا همه مقادیر یکسانی خرید کنند، این گروه‌بندی مفید نیست. خوشبندی هدف به دنبال یافتن خوشه‌هایی است که با توجه به مقدار یک متغیر هدف خاص متفاوت هستند. این کار با افزایش وزن (اهمیت) متغیر هدف در معیار فاصله یا کُد کردن متغیرها با توجه به هدف انجام می‌شود. در حالت حدی که فقط متغیر هدف در معیار فاصله در نظر گرفته می‌شود، خوشبندی شبیه دسته‌بندی می‌شود. فرق مهم خوشبندی هدف و دسته‌بندی، این است که مقوله از پیش تعریف شده‌ای برای خوشبندی هدف وجود ندارد.

### رگرسیون و سریهای زمانی (در مدیریت ارتباط با مشتری)

روش رگرسیون سعی در پیش‌بینی یک خروجی پیوسته با استفاده از یک تابع دارد که میزان خطا از الگو را در داده‌ها نشان می‌دهد. روش‌های رگرسیون را می‌توان برای کاربرد دسته‌بندی دوتایی و همچنین برای امتیازدهی و پیش‌بینی استفاده کرد. روش‌های سری زمانی مثل رگرسیون مقادیر پیوسته را پیش‌بینی می‌کنند، ولی در طول زمان، گرایشها و چرخه رفتار را نیز مدل می‌کنند.

### قواعد تلازمی (در مدیریت ارتباط با مشتری)

قواعد تلازمی و الگوهای مکرر، وقایعی که در مجموعه داده اتفاق می‌افتد را توصیف می‌کنند. قواعد تلازمی، ارتباط بین موارد موجود در یک مجموعه داده است بدون اینکه ترتیب زمانی یا ترتیب خاصی داشته باشند.

## فیلتر کردن مشارکتی<sup>۱</sup>

فیلتر کردن مشارکتی وقایع را برای یافتن مجموعه‌هایی که شبیه یکدیگرند تحلیل می‌کنند. در مفهوم CRM، وقایع معمولاً بر اساس مشتریان گروه‌بندی می‌شود. ایده اصلی این است که اگر مشتریان با سابقه خرید مشابه کالای خاصی را بخزند، ممکن است مشتری جدید دارای همان مشخصات نیز همان کالا را بخرد. برای هر مشتری در مجموعه داده، با روش فیلتر کردن مشارکتی یک گروه از مشتریان پیدا می‌شوند که گزارشات وقایع‌شان به یکدیگر شبیه هستند. وقتی گروه شکل می‌گیرد، وقایع می‌تواند با چیزی که در کل گروه عمومیت دارد، امتیازدهی شود. خروجی فیلتر کردن مشارکتی می‌تواند هم مشتریان شبیه به هم و هم وقایع امتیاز داده شده باشد. برای مثال، اگر واقعه، خرید محصولات باشد، امتیازدهی وقایع می‌تواند برای توصیه محصولات به کار بrede شود. هر محصولی که بیشترین امتیاز را بیاورد توصیه می‌شود. فیلتر کردن مشارکتی مشابه خوشه‌بندی برای کشف گروه‌های طبیعی در مجموعه داده است. تفاوت مهم فیلتر کردن مشارکتی این است که گروه‌های متفاوتی گردآگرد یکدیگر در مجموعه داده شکل می‌گیرند. در واقع هر مشتری مرکز یک گروه واحد است و مشتری A ممکن است در گروه مشتری B ظاهر شود ولی عکس آن درست نباشد. در نتیجه روش فیلتر کردن مشارکتی برای شخصی‌سازی سیستم به کار بrede می‌شود.

## منابع

- 1) Berson A. , Smith S. ,Thearling K. (2001) "Bulding Data Mining Application for CRM" , McGraw-Hill.
  - 2) Ye N. (2003) "THE HANDBOOK OF DATA MINING", LAWRENCE ERLBAUM ASSOCIATES , PUBLISHERS Mahwah, New Jersey London.
  - 3) Rygielski C. , Wang J. C. , Yen D. C. (2002), "Data mining techniques for customer relationship management", *Technology in Society* 24 (2002) 483–502.
  - 4) Berry M. J. A. , Linoff G. S. (2004) *Data Mining Techniques for Marketing, Sales, and Cutomer Relationship Management*, Second edition, Wiley.
- ۵) سعیدرضا ملک محمدی (۱۳۸۵)، سمینار کارشناسی ارشد، کاربرد داده‌کاوی در مدیریت ارتباط با مشتری  
دانشگاه علم و صنعت، استاد راهنمای دکتر مهدی غضنفری.

## واژه نامه

Run tests	آزمونهای ردیف
Dirty	آشagnetگی و آلودگی
Deviance Statistic	آماره انحراف
Entropy	آنتروپی
Grass- Root Political Initiative	ابتکار سیاسی اجتماعی محلی
Hypercube	آئر مکعب
Hypertext	ا!ر متنها
Lateral	اتصال جانبی
Community Mining	اجتماع کاوی
On Line Forums	اجتماعات برخط
Web Communities	اجتماعات وبی
Record Linkage	ارتباط رکورد
Topologic	ارتباط مکانی
Network Values	ارزش شبکه‌ای
Life Time Value	LTV ارزش طول عمر
Extract	استخراج
ETL (Extraction/Transformation>Loading)	استخراج، تبدیل، بارگذاری
Case Base Reasoning	استدلال مبتنی بر مورد
Customer Focused Strategy	استراتژی تمرکز بر مشتری
Model Induction	استقراء مدل
Decision Tree Induction	استنتاج درخت تصمیم
Nominal	اسمی
Objects	اشیاء
Information	اطلاعات
Confidence	اطمینان
Noise	اغتشاش

Redundancy	افزونگی
Redundant	افزونه
Marketing Action	اقدام بازاریابی
Meningitis	التهاب مغزی
Preferential Attachment Model	الحاق ترجیحی
Patterns	الگوهای
Frequent Pattern	الگوهای مکرر
Nodes Degrees	امتیاز گره‌ها
Enterprise Warehouse	انباره بنگاه
Data Warehouse	انباره داده‌ها
Virtual Warehouse	انباره مجازی
Feature Selection	انتخاب مشخصه‌ها
Relation Selection & Extraction	انتخاب و استخراج رابطه
Transform	انتقال
Deviation	انحراف
Notions	ایده‌ها

**ب**

Load	بارگذاری
Data Mart	بازارچه داده‌ها
Mass Marketing	بازاریابی انبوه
Direct Marketing	بازاریابی مستقیم
Viral Marketing	بازاریابی ویروسی
Resubstitution	بازجانشانی
Information Retrieval	بازیافت اطلاعات
Data Archaeology	پاسدان شناسی داده‌ها
Segment	بخش
Segmentation	بخش بندی
Lossless	بدون اتلاف
Novel	بدیع

Novelty	بدیع بودن
Fitness	برازندگی
Label	برچسب
Class Label	برچسب دسته
Inductive Logic Programming	برنامه‌ریزی منطقی قیاسی
Longest Common Subsequences Similarity	
LCSS	بزرگترین زیردنباله‌های مشترک
Gridlock	بن بست
Well Understood	به خوبی قابل درک
BMU: Best Matching Unit	بهترین واحد انطباق
Refresh	بهروزرسانی
Timely	بهموقع
Global Optimum	بهینه سراسری
Local Optimum	بهینه محلی
Non Stationary	محی ثبات
Over Fit	بیش برداش

## پ

Data Cleaning	پاکسازی داده ها
NON Volatile	پایابی
Robustness	پایداری
Dispersion	پراکندگی
On_line Analytical Processing	بردازش تحلیلی برخط
Profile	پروفایل
Post-Process	پس پردازش
Back Propagation	پس انتشار
Support	پشتیبان
Time Variant	پویاپذیری
Preprocess	بیش برداش

Prediction	پیش‌بینی
Centroid	پیوند مرکزی
Instant Messaging	پام آنی
Predictive	پیش‌بینی
Random Predictor	پیشگویی تصادفی
Linkage	پیوند
Single Link	پیوند تکی
Complete Link	پیوند کامل
Average Link	پیوند متوسط
Link Mining	پیوند کاوی
Out-Links	پیوندهای خروجی

## ت

Dynamic Time Warping: DTW	تاباندن زمانی پویا
Word of Mouth	تأثیر مثبت گفتار شفاهی
Transformation	تبديل
Data Transformations	تبديلات داده
Hill Climbing	تپه توردي
Consolidation	تبسيط
Discrete Fourier Transform: DFT	تجزیه فوریه
Singular Value Decomposition: SVD	تجزیه مقدار منفرد
Discrete Wavelet Transform (DWT)	تجزیه موجک
Divisive	تجزیه‌ای یا تقسیمی
Aggregation	تجمیع
Feature Aggregation	تجمیع مشخصه‌ها
Agglomerative	تجمیعی
Market Basket -Basket Data	تحلیل سبد بازار
Principal Component Analysis: PCA	تحلیل مولفه‌های اصلی
Line Segment Approximation	تخمین قطعه‌ای خط
Piecewise Linear Approximations: PLA	تخمینهای خطی قطعه‌ای

Lexicographic Order	ترتيب حروف الفبا
Combination	ترکیب
Pattern Recognition	تشخیص الگو
Discrepancy Detection	تشخیص مغایرت
Reformulation Schema	تشکیل مجدد شماتیک
Decision Science	تصمیم‌گیری علمی
Data Projection	تصویر کردن
Global Representation	تصویر کلی
Snapshot	تصویر آنی
Projections	تصویر کردن
Random Projection	تصویر کردن تصادفی
Conflicts	تضادها
Adaptation	تطبیق
Generalization	تممیم
Interpretability	تفسیر پذیری
Descriptive Data Summarization	تلخیص توصیفی داده‌ها
Lazy	تبیل
Activation Function	توابع فعال سازی
Heavy-Tailed	توزیع‌های دمپهن
Descriptive	توصیفی
Attribute Value Description	توضیحات ویژگی - ارزش

## ث

ثروتمند، ثروتمندتر می‌شود

## ج

Density Attractors	جادذب چگالی
Interestingness	جالب بودن
Exhaustive	جامع
Integrated	جامعیت

Conditional Probability Table	جدول احتمال شرطی
Contingency Table	جدول تصادفی
Biomass	جرم حیاتی
Clickstream	حریان کلیکها
Spurious	جملی
Shrinking Diameter	جمع شدن قطر
Compactness	جمع و جور بودن مدل

**ج**

Quartiles	چارکها
Quantile	چندک
Quantile –Quantile (Q-Q)	چندک چندک

**ح**

Sensitivity	حساسیت
Facts	حقایق

**خ**

Deliberate Errors	خطاهای عمدی
Summarization	خلاصه سازی
Clustering	خوشبندی

Density- Based Spatial Clustering of Applications with Noise	خوشبندی فضایی بر پایه چگالی برای داده های معشوش
---	--

**د**

Data	داده
Neural Network Data Mining	داده کاوی بر مبنای شبکه های عصبی مصنوعی
Outliers	داده های پرت

Data Scrubbing	داده رویی
Missing Data	داده های مفقوده
Sufficient Data	داده کافی
Interquartile range (IQR)	دامنه میان چارکی
Knowledge	دانش
Degree of Interest	درجه جذبیت
Data Harvesting	درو کردن داده ها
Classification	دسته بندی
Predefined Classes	دسته های از پیش تعیین شده
Discrete Classes	دسته های گسسته
Precision	دقت
Accurate	دقیق
Small World	دبیای کوچک

## ر

Ordinal	رتبه ای
Competition	رقابت
Regression	رگرسیون
Binning	روش بسته بندی
A Statistical Information Grid Approach	روش شبکه اطلاعات آماری
FastMap	روش نگاشت سریع
Split Selection	روشهای انتخاب نقطه انشعاب
Hierarchical	روشهای سلسله مراتبی
Density based	روشهای مبنی بر چگالی
Impurity – Based	روشهای مبنی بر ناخالصی

## ز

Standard Query Language	زبان پرس و چوی استاندارد
Time- Lagged	زمان تأخیری
Subgraph	زیر گرافها

**س**

Feature Construction	ساخت مشخصه ها
Attribute Construction	ساخت ویژگی
Structural	ساختاری
Simplicity	سادگی
Consistent	سازگار
Knowledge- Sharing Sites	سایتهای شرکت دانش
Field Overloading	سر بر ا شدن فیلد
Speed	سرعت
Bucket	سطله اها
Information Gain	سود اطلاعاتی
Biased	سوگیری
Data Base Management System	سیستم مدیریت پایگاه داده
Management Information System (MIS)	سیستمهای اطلاعات مدیریت
Complex Systems	سیستمهای پیچیده
Legacy	سیستمهای عملیاتی یا میراثی
Source System	سیستمهای منبع
Sigmoid Scaling	سیگموئید

**ش**

Gini Index	شاخص جینی
Probabilistic Dependency Network	شبکه های وابستگی احتمالی
Artificial Neural Networks	شبکه های عصبی مصنوعی
Multi Relational Social Network	شبکه اجتماعی چند رابطه ای
Web of Trust	شبکه اعتماد
Feed Forward	شبکه چند لایه پیشخور
Competitive Network	شبکه رقابتی
Lattice	شبکه نرده بان
Six Degrees of Separation	شش سطح جدایی

Specificity	شفافیت
Enumeration	شمارش

	ص
Accuracy	صحت

	ط
Categorical	طبقه ای

	ع
Equal-width	عرض ثابت
Equal-depth	عمق ثابت
Pruning	عملیات هرس

	ف
Reach ability- Distance	فاصله دسترسی
Average Distance	فاصله متوسط
Core- Distance	فاصله مرکزی
Interval	فاصله ای
Usefulness	فایده
Iterative	فرایندی تکراری
Closed Versus Open World Assumption	فرض دنیای باز در مقابل دنیای بسته
Up-Selling	فروش بالا سری
Cross-Selling	فروش کناری
Campaign	فعالیت تبلیغی
Superlinearly	فرق خطی
Data Understanding	فهم داده
Business Understanding	فهم کسب و کار
Collaborative Filtering	فیلتر کردن مشارکتی

**ق**

Density Reachable	قابل دسترس چگال
Densification Power Law	قانون تراکم توانی
Null Rule	قانون نهی
Growth Power Law	قانون رشد توانی
Unique Rule	قانون یکتاپی
Effective Diameter	فطر مؤثر
Induced Rules	قواعد استنتاج شده
Association Rules	قواعد تلازمی
Strong Rule	قواعد قوی
Frequent Rule	قواعد مکرر

**ک**

CART	کارت
Fully Connected	کاملاً متصل
Reduction	کاهش
Dimensionality Reduction	کاهش بعد نقاط
Numerical Reduction	کاهش تعداد نقاط
Knowledge Acquisition	کسب دانش
Knowledge Discovery and Data Mining (KDD)	کشف دانش و داده‌کاوی
Hyperlink	كلمات ابزیوندی
Quantitative	كمی
Quality-of-Life	کیفیت زندگی
Bibliography	كتاب‌شناسی
Vector Quantization	كمی‌سازی برداری
Coupon	کوبن

## Line Graph

گ

Acyclic Graph	گراف خطی
Central Tendency	گرایش مرکزی
Use Net Groups	گروههای خبری
Ambassador	گره سفیر
Discrimination	گستته سازی
Evidence	گواهی
Diversity	گوناگونی

ل

Data Dredging	لایروبی داده ها
Wrapper	لغاف

م

Confusion Matrix	ماتریس اختشاش
Support Vector Machine (SVM)	ماشینهای بردار پشتیبان

## Cross Industry Standard Process for Data

متداول‌زی CRISP

Mining	متصل چگال
Density Connected	متقابل ناسازگار
Exclusive	متن لنگر
Anchor Text	مثبت درست
True Positive	مثبت غلط
True Negative	مجموع تمام مسیرها
Ensemble of All Paths	مجموع مربع خطاهای ناهمگن
Sum Square Error	مجموعه داده‌های کاهش یافته
Heterogeneous	محدودیت‌های شبیب
Reduced Set	
Slope Constraints	

Metadata Repository	مخزن فراداده
Central Repository	مخزن مرکزی
Mode	مد
Water Fall	مدل آبشاری
Memory Based Reasoning	مدل استدلال بر مبنای حافظه
Snow-Flake Schema	مدل برفدانه
Response Model	مدل پاسخ
Star Schema	مدل ستاره‌ای
Star Net	مدل شبکه ستاره‌ای
Fact Constellation	مدل صورت فلکی
Copying Model	مدل کپی کننده
Spiral	مدل مارپیچی
Analytical CRM	مدیریت ارتباط با مشتری تحلیلی
Operational CRM	مدیریت ارتباط با مشتری عملیاتی
Collaborative CRM	مدیریت ارتباط با مشتری مشارکتی
Customer Relationship Management (CRM)	مدیریت ارتباط با مشتری
Ordering Points To Identify the Clustering Structure	مرتب‌سازی نقاط برای شناسایی ساختار خوشه‌بندی
Generalization Problem	مسئله تعمیم
Directly Density Reachable	مستقیماً قابل دسترس چگال
Eager	مشتاق
Customer Driven	مشتری محوری
Customer Centric	مشتری مداری
Prospects Customers	مشتریان احتمالی
Loyal Customer	مشتریان وفادار
Feature	مشخصه
Reconciliation	مصالحه
Data Decay	صرف گذشته
Visualization	تصویرسازی
Data Visualization	تصویرسازی داده‌ها

The Curse Of Dimensionality	مسئیت بعد
Significant	معنی دار
Noisy	مغشوش
Missing Value	مقادیر مفقوده
Normalization by decimal scaling	مقیاس بندی اعشاری
Multidimensional Scaling: MDS	مقیاس بندی چند بعدی
Data Cube	مکعب داده
Data Auditing	ممیزی داده
True Negative	منفی درست
False Negative	منفی غلط
Subject Oriented	موضوع محور
Midrange	میان دامنه
Interdisciplinary	میان رشته ای
Mean	میانگین
Median	میانه
Constant Average Degree Assumption	میانگین امتیاز ثابت

## ن

Inconsistent	ناسازگار
Incomplete	ناقص
Syntactical	نحوی
Normalization	نرمال سازی
Neuron	نرون
K Nearest Neighborhood	نزدیک ترین همسایگی
Ratio	نسبتی یا نسبی
Minpts	نقاط
Poverty Map	نقشه فقر
Self-Organizing Maps: SOM	نقشه های خودسازمان یا خودسازمانده
Kohonen Maps	نقشه های کوهونن
SOFM: Self-Organizing Feature Maps	نقشه های مشخصه خودسازمان

Schema Mapping	نگاشت شماتیک
Symbolic	نمادین
Scatter Plot	نمودار پراکش
Quantile Plot (Q-P)	نمودار چندک
Dendogram	نمودار درختانهای
Bar Chart	نمودار ستونی
Pie Chart	نمودار کلوچه‌ای
Loess Curve	نمودار لوئس
Sampling	نمونه برداری
Instances	نمونه‌ها

## و

Dependency	وابستگی
Probabilistic Dependencies	وابستگی‌های احتمالی
Garbage in Garbage Out	ورودی نامناسب، خروجی نامناسب
Connection Weight	وزن اتصال
Slope Weighting	وزن‌دهی به شبیب

## ه

Spam	هرزنگاری
Prune	هرس
Prepruning	هرس اولیه
Postpruning	هرس ثانویه
Ontologies	هستیان‌شناسی
Correlation	همبستگیها
Overlapping	همپوشانی
Common Neighbors	همسایگان مشترک
Cooperation	همکاری
Smooth	هموار
Smoothing	هموارسازی

## ی

Instance Based Learner	یادگیر مبتنی بر نمونه
Machine Learning	یادگیری ماشین
Work File	یک فایل آماده برای کار
Relational Learning	یادگیری رابطه‌ای بالها
Edges	
Center Defined Cluster	یک خوشه مرکز محور
Collective Consolidation	یکسازی جمعی

