

WILEY TIMELY. PRACTICAL. RELIABLE.

Data Mining Techniques

Second Edition

For Marketing, Sales, and Customer Relationship Management

Michael J. A. Berry
Gordon S. Linoff

Jiawei Han
Micheline Kamber

Data Mining

Concepts and Techniques

داده کاوی

سمیه علیزاده

مطالب مورد بحث در کلاس

- تعاریف و مفاهیم
- انباره داده ها
- آماده سازی داده ها (پیش پردازش داده ها)

مطالب مورد بحث در کلاس

- خوشه بندی
- دسته بندی
- قوانین انجمنی
- سریهای زمانی
- وب کاوی
- متن کاوی
- پیوندکاوی و تحلیل شبکه های اجتماعی

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مطالب مورد بحث در کلاس

- متدولوژی اجرای پروژه های داده کاوی
- کاربردهای داده کاوی در بازاریابی
- کاربردهای داده کاوی در مدیریت ارتباط با مشتری
- امنیت در داده کاوی

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مبانی انباره داده ها



- ◆ تاریخچه و تعاریف و مفاهیم
- ◆ ویژگیهای انباره داده ها
- ◆ ساختار انباره داده ها

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مبانی انباره داده ها



- ◆ تاریخچه و تعاریف و مفاهیم
- ◆ ویژگیهای انباره داده ها
- ◆ ساختار انباره داده ها

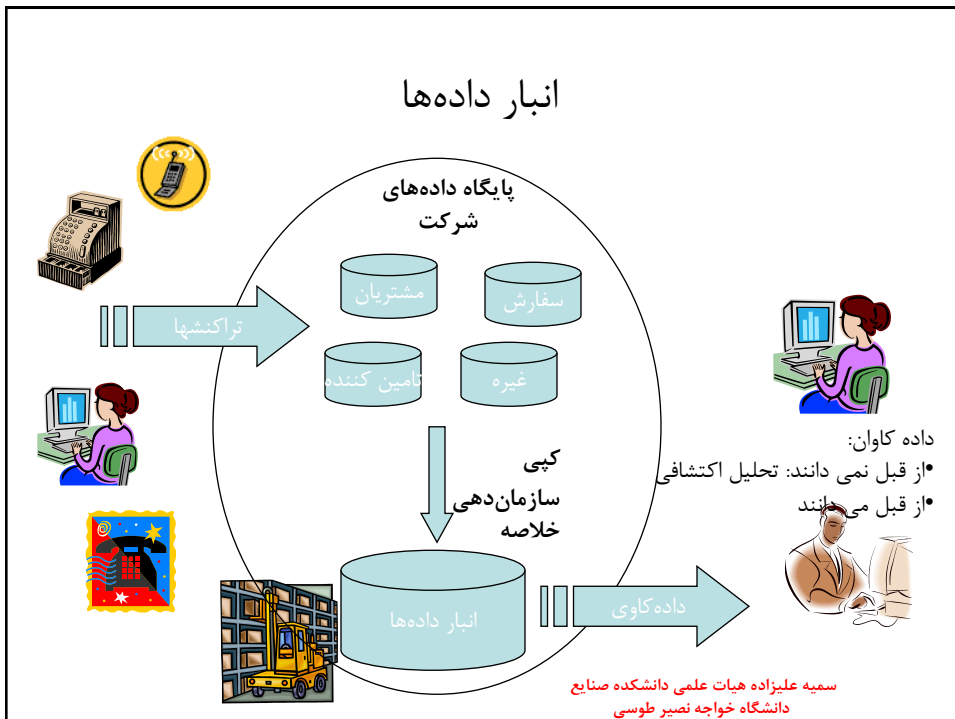
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

تاریخچه

- در دهه ۹۰ میلادی پدیده انبار داده ها ظهور یافت . قبل از انبار سازی داده ها سیستم های کامپیوتری جهت ذخیره ، جمع آوری ، تغییر و تصحیح بیت های داده ها طراحی شده بودند. این سیستم های اولیه به سیستم های عملیاتی یا میراثی موسوم هستند .

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

انبار داده ها



انبار داده یک ضرورت کامل برای داده کاوی مؤثر می باشد .

- انبار داده ها ، داده های خام را برای تحلیل به گونه ای بهینه آماده می سازند . این آماده سازی به شکلهای مختلف سودمند می باشد .
- یکی از ماهیت های وجودی انبار داده ها ، یکپارچگی داده ها هنگام قرار گیری در انبار داده هاست . این بدین معنی است که دقت بسیاری بکار گرفته می شود تا یکنواختی و پیوستگی در درک اهداف عام سازمانی بوجود آید .
- اگر انبار داده ها وجود نداشته باشد ، داده کاو بایستی زمان بسیار زیادی را برای جمع آوری ، پاکسازی ، تمیز کردن و یکپارچه سازی داده ها سپری کند . بدین ترتیب وقت بسیاری باید صرف کرد تا بتوان کار تحلیل داده ها را آغاز نمود .

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

انبار داده یک ضرورت کامل برای داده کاوی مؤثر می باشد .

- دومین علتی که می توان برای موفقیت داده کاوی به خاطر وجود انبار داده ها ذکر کرد این است که در انبار داده ها ، داده های تاریخی جمع آوری و سازماندهی می شوند .
- وجود داده های تاریخی برای یافتن الگوها و روابطی که سازمان بدنبال آنهاست ، برای داده کاوی یک ضرورت است .
- اگر چنانچه این داده های تاریخی وجود نداشته باشند ، داده کاو بایستی بدنبال جمع آوری آنها باشد .

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

انبار داده یک ضرورت کامل برای داده کاوی مؤثر می باشد .

- علت سوم اهمیت انبار داده ها برای داده کاوی مؤثر، این است که انبار داده ها شامل داده های جزئی و داده های کلی، در کنار یکدیگر می باشد .
- بدون تردید ، داده کاو به اطلاعات جزئی برای تحلیل نیازمند است، اما داده های خلاصه شده نیز بکار می آیند .
- هنگامی که نمونه ای از انواع داده های خلاصه شده وجود داشته باشد ، داده کاو می تواند به سرعت بررسی کند که چه چیزی در انبار داده ها هست ، یا چیزی نیست و این باعث کاهش تکرار تحلیل ها توسط داده کاو می شود .

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مبانی انبار داده ها



- ◆ تاریخچه و تعاریف و مفاهیم
- ◆ ویژگیهای انبار داده ها
- ◆ ساختار انبار داده ها

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

ویژگیهای انباره داده ها

- موضوع محوری (Subject Oriented)
- جامعیت (Integrated)
- مهم بودن عامل زمان (Time Variant)
- غیر فرار و دایمی بودن (NON Volatile)

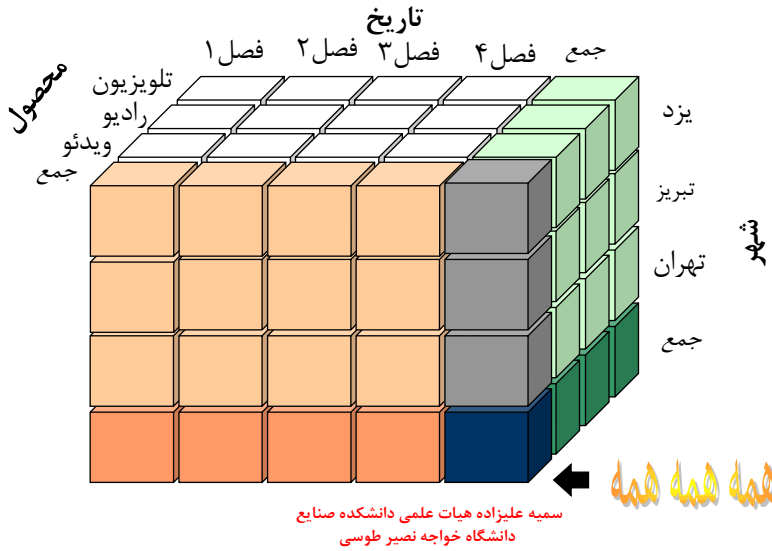
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

موضوع محوری

- داده ها طبق یک موضوع خاصی سازماندهی میشوند، به عنوان مثال داده های مربوط به مشتریان ، محصولات و یا داده های مرتبط با فروش هر کدام جداگانه در نظر گرفته میشوند. اما در پایگاه داده های معمولی، داده ها بر اساس عملیات و پردازش های روزانه ایجاد میگردند و موضوع آنها مرتبط با کل پردازش میباشد.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

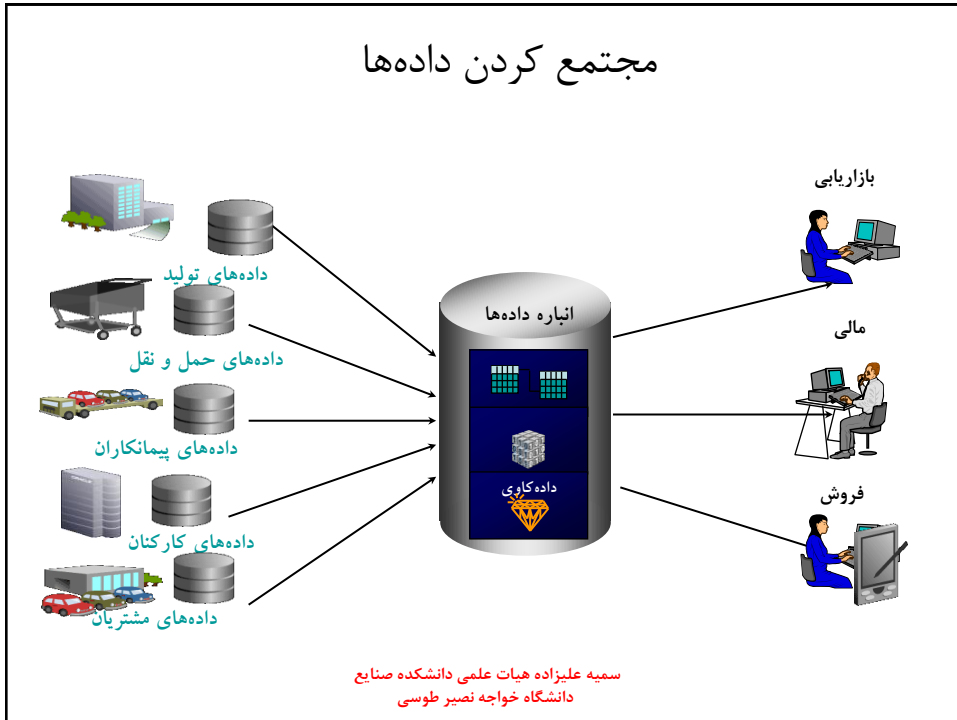
خلاصه کردن داده‌ها



جامعیت

- داده های انباره های داده ای از تجمع دیگر داده ها ساخته میشوند. این داده ها ممکن است مربوط به پایگاه داده های رابطه ای، فایل های بدون ساختار و یا رکوردهای مرتبط با پردازش های Online باشد.

مجموع کردن داده‌ها

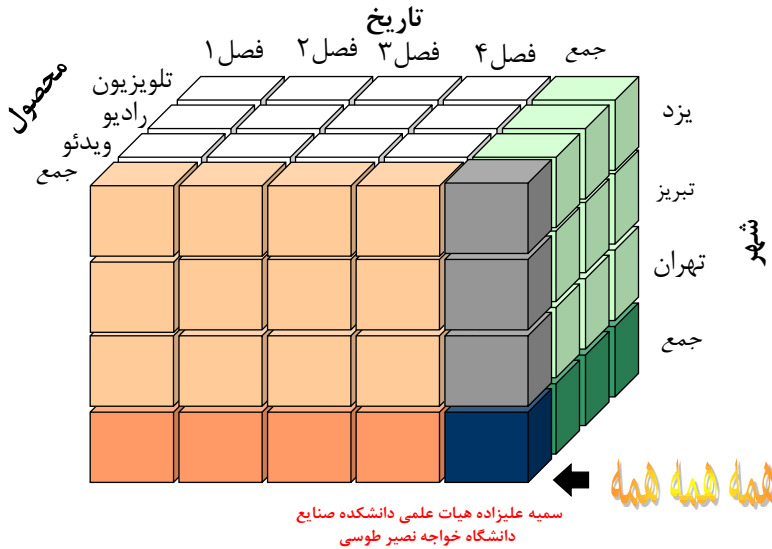


عامل زمان

- افق زمانی برای انبار داده‌ها بسیار مهمتر از داده‌های مرتبط با سیستم‌های عملیاتی میباشد. در ساختار انبار داده‌ها عاملی به نام زمان در نظر گرفته میشود. این عامل زمانی میتواند بصورت یک عامل ضمنی و یا به وضوح بیان گردد. اما در سیستم‌های عملیاتی عامل زمان عاملی کلیدی نیست.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

خلاصه کردن داده‌ها



غیر فرار بودن

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

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

منافع اصلی انباره داده ها

- منافع اصلی از قابلیت های انباره داده عبارتند از :
 - به سرعت قابل دسترسی هستند ، زیرا از لحاظ فیزیکی در یک مکان قرار دارند .
 - بوسیله کاربر نهایی به راحتی قابل حصول هستند .

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مبانی انباره داده ها



- ◆ تاریخچه و تعاریف و مفاهیم
- ◆ ویژگیهای انباره داده ها
- ◆ ساختار و معماری انباره داده ها

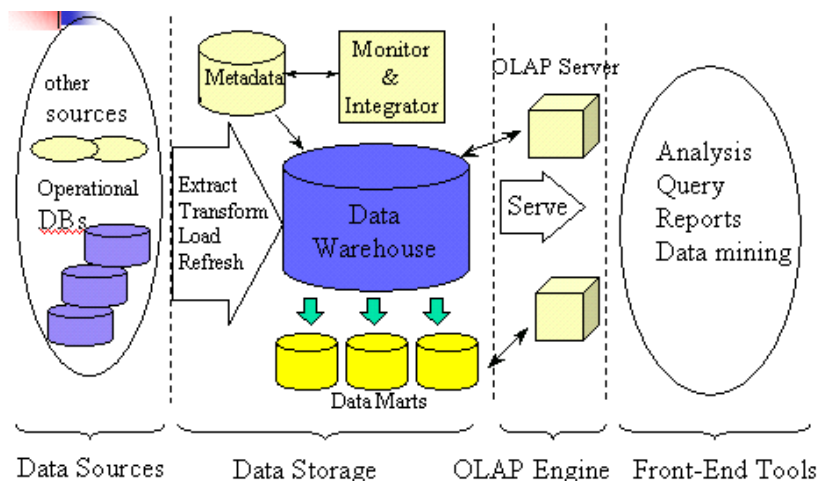
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

معماری انبار داده

- جهت تهیه و کار با انبار داده ها یک ساختار چند لایه ای وجود دارد که در شکل بعد نشان داده شده است. در ابتدا با چهار عمل اصلی Extract, Transform, Load, Refresh داده ها از پایگاه داده های معمولی (Data Source) که از پردازش های عملیاتی سیستمها ایجاد شده اند، جمع آوری شده و به انبار داده فرستاده میشود.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

معماری انبار داده



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

تعاریف

- داده های موجود در انباره داده ها مستقیماً با دیگر سیستمها در ارتباط نیستند و اگر یک سیستم عملیاتی بخواهد از داده های انباره داده استفاده کند از **Data Mart** ها که به طور موقت شامل بخشی از داده ها را که مرتبط با آن نوع سیستم خاص میباشند، استفاده میکند.
- متاداده ها ، داده های مربوط به داده ها هستند . متاداده ها وضعیت های مختلف داده ها را توصیف می کنند

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

تعریف (Data Mart)

- انباره داده کوچک محدود شده در ابعاد
Mini-warehouses, limited in scope
– یک **Data Mart** یک زیر مجموعه کوچک و محدود از یک **DW** می باشد که دارای مجموعه ای از موضوعات برای یک بخش از کاربران است.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

پردازشها در انبار داده ها

- انبار داده ها نیازی به پردازشهایی از قبیل : تراکنشهای داده ای، بازیافت و مکانیزمهای کنترل همزمان ندارد. تنها اعمالی که در انبار داده ها صورت میپذیرد عبارتند از: بار گذاری (مقدار دهی) اولیه داده ها و دسترسی به داده ها.
- پردازشی که بر روی داده های انبار داده ها انجام میگیرد OLAP نامیده میشود (OLAP در مقابل OLTP که پردازش هایی است که بر روی داده های پایگاه داده ها انجام میگیرد آمده است). جدول بعد OLAP و OLTP را با توجه به برخی پارامترهای مهم مقایسه کرده است:

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مقایسه پردازشها در داده کاوی با داده های معمولی

| ویژگیها | OLTP | OLAP |
|--------------------|-----------------------------|-------------------------|
| کاربران | اپراتورها | کارشناسان خبره |
| کارکرد | کارهای روزمره | تصمیم گیری |
| طراحی پایگاه داده | بر مبنای کاربرد | بر مبنای موضوع |
| داده | جاری، روزبه-روز و با جزئیات | تاریخی، چند بعدی، مجتمع |
| کاربرد | تکراری | در موارد خاص |
| نحوه دسترسی | خواندن و نوشتن | کاوش و کشف |
| واحد کاری | پردازش ساده | جستجوهای پیچیده |
| تعداد رکوردها | دهها | میلیونها |
| تعداد کاربران | هزاران | صدها |
| اندازه پایگاه داده | MB_GB ۱۰۰ | TB_GB ۱۰۰ |
| شاخص | پردازش | جستجو و پاسخ |

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

ایجاد و نگهداری یک انباره داده ها

انباره های داده به چندین ابزار نیاز دارند تا وظایف ذیل را خودکار کرده یا مورد پشتیبانی قرار دهند:

- استخراج داده ها (Data extraction) از منابع داده متفاوت خارجی
- پالایش داده ها (Data cleaning): پیدا کردن و حل ناسازگاریها از داده های منبع
- یکپارچگی و تبدیل داده ها: بین فرمتها و زبانهای مختلف

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

ایجاد و نگهداری یک انباره داده ها

- بارگیری داده ها (Data loading): بارگیری داده ها در انباره داده ها
- تکرار داده ها (Data replication): از پایگاه داده منبع به انباره داده ها
- تازه کردن داده ها (Data refreshment)
- آرشیو داده ها (Data archiving)
- چک کردن کیفیت داده ها
- تحلیل متا داده ها (metadata)

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مدلسازی مفهومی انباره داده ها

سه شمای مفهومی پایه عبارتند از :

- شمای ستاره (**Star schema**)
- شمای دانه برفی (**Snowflake schema**)
- شمای فلکی (**Fact constellations**)

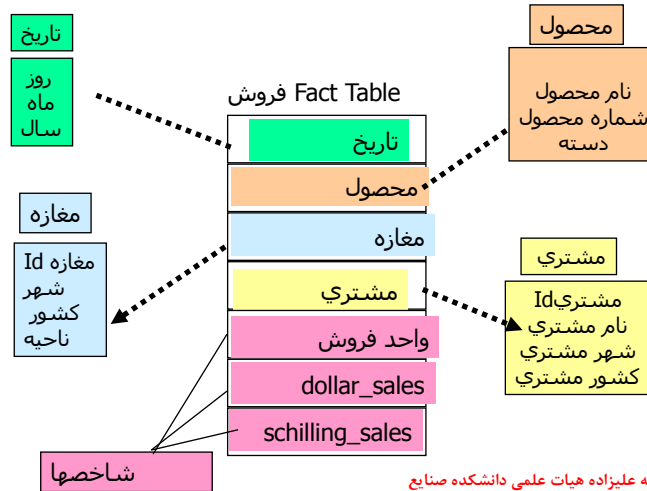
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

شمای ستاره

• شمای ستاره: یک **fact table** در مرکز به تعدادی از
فهرستهای ابعادی (**dimension tables**) متصل می
شود .

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

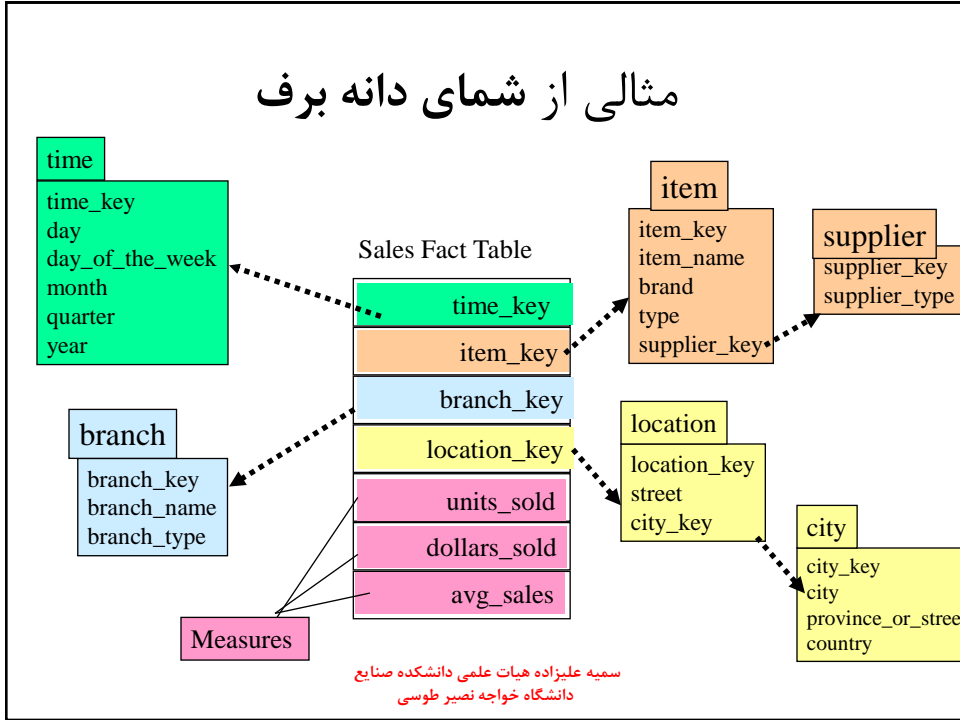
مثالی از شمای ستاره



شمای دانه برف

شمای دانه برف: یک پالایشی از شمای ستاره است که سلسله مراتب ابعاد بصورت صریح با نرمال کردن فهرستهای ابعادی (dimension tables) نمایش داده می شود.

مثالی از شمای دانه برف



مثالی از شمای صورت فلکی

