

# **Data Analysis and OLAP**

Chapter 18: Data Analysis and Mining  
Database system concepts 5th Edition  
Silberschatz, Korth and Sudarshan

# Data Analysis and OLAP

## ■ Online Analytical Processing (OLAP)

- Interactive analysis of data, allowing data to be summarized and viewed in different ways in an online fashion (with negligible delay)

## ■ Data that can be modeled as dimension attributes and measure attributes are called **multidimensional data**.

### ● **Measure attributes**

- ▶ measure some value
- ▶ can be aggregated upon
- ▶ e.g. the attribute *number* of the *sales* relation

### ● **Dimension attributes**

- ▶ define the dimensions on which measure attributes (or aggregates thereof) are viewed
- ▶ e.g. the attributes *item\_name*, *color*, and *size* of the *sales* relation

# Cross Tabulation of sales by *item-name* and *color*

size:

*color*

	dark	pastel	white	Total
<i>item-name</i>				
skirt	8	35	10	53
dress	20	10	5	35
shirt	14	7	28	49
pant	20	2	5	27
Total	62	54	48	164

- The table above is an example of a **cross-tabulation** (**cross-tab**), also referred to as a **pivot-table**.
  - Values for one of the dimension attributes form the row headers
  - Values for another dimension attribute form the column headers
  - Other dimension attributes are listed on top
  - Values in individual cells are (aggregates of) the values of the dimension attributes that specify the cell.

# Relational Representation of Cross-tabs

- Cross-tabs can be represented as relations
  - We use the value **all** is used to represent aggregates
  - The SQL:1999 standard actually uses null values in place of **all** despite confusion with regular null values

<i>item-name</i>	<i>color</i>	<i>number</i>
skirt	dark	8
skirt	pastel	35
skirt	white	10
skirt	<b>all</b>	53
dress	dark	20
dress	pastel	10
dress	white	5
dress	<b>all</b>	35
shirt	dark	14
shirt	pastel	7
shirt	white	28
shirt	<b>all</b>	49
pant	dark	20
pant	pastel	2
pant	white	5
pant	<b>all</b>	27
<b>all</b>	dark	62
<b>all</b>	pastel	54
<b>all</b>	white	48
<b>all</b>	<b>all</b>	164

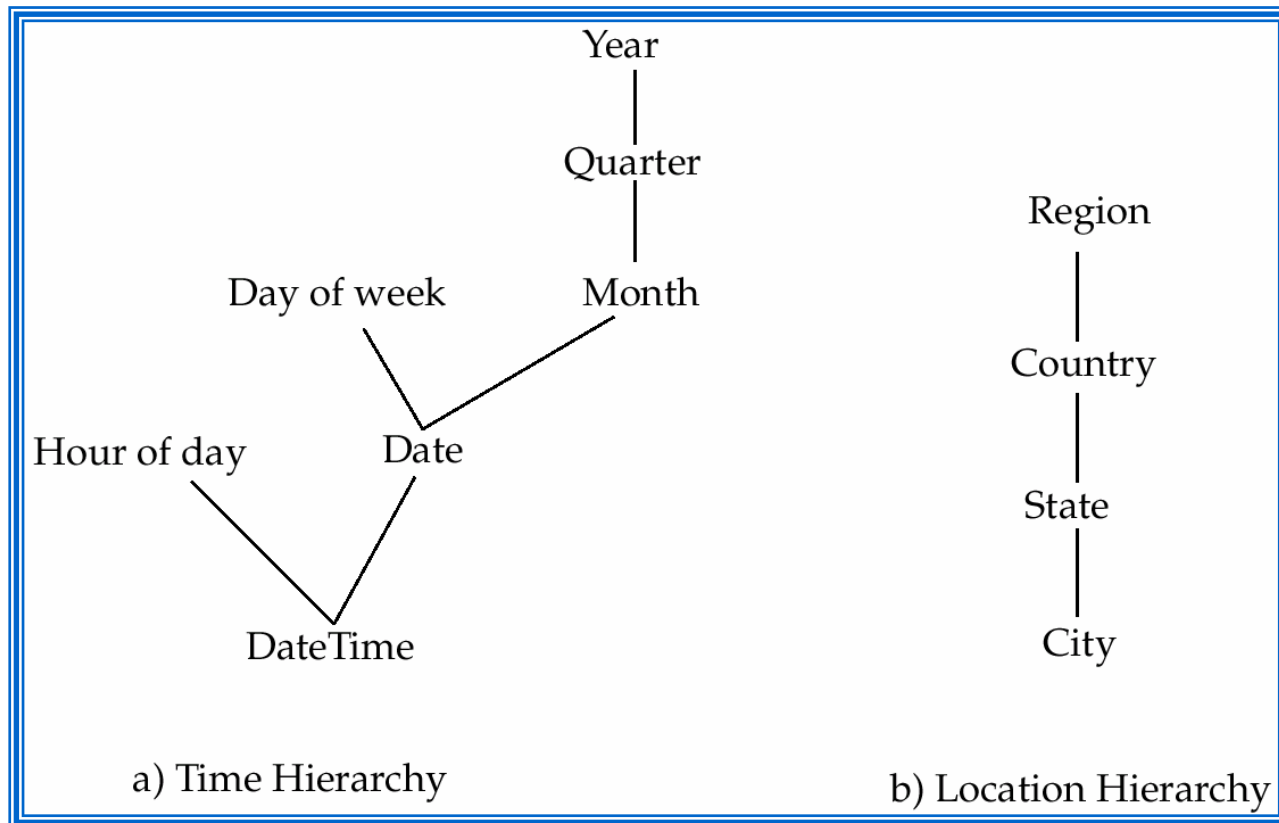
# Data Cube

- A **data cube** is a multidimensional generalization of a cross-tab
- Can have  $n$  dimensions; we show 3 below
- Cross-tabs can be used as views on a data cube

color	item name					size		
	skirt	dress	shirts	pant	all	small	medium	large
dark	8	20	14	20	62	4	18	16
pastel	35	10	7	2	54	9	45	34
white	10	8	28	5	48	42	21	77
all	53	35	49	27	164	55	84	111

# Hierarchies on Dimensions

- **Hierarchy** on dimension attributes: lets dimensions to be viewed at different levels of detail
  - 👉 E.g. the dimension DateTime can be used to aggregate by hour of day, date, day of week, month, quarter or year



# Cross Tabulation With Hierarchy

- Cross-tabs can be easily extended to deal with hierarchies
  - ☞ Can drill down or roll up on a hierarchy

<i>category</i>	<i>item-name</i>	dark	pastel	white	total	
womenswear	skirt	8	8	10	53	
	dress	20	20	5	35	
	subtotal	28	28	15		88
menswear	pants	14	14	28	49	
	shirt	20	20	5	27	
	subtotal	34	34	33		76
total		62	62	48		164