# CSX4202/ITX4202: Data Mining

# Lecture 3

Asst. Prof. Dr. Rachsuda Setthawong

Computer Science Department

Assumption University

# Outlines

- Data Exploration: Tasks and Techniques
  - Summary Statistics
  - Data Visualization
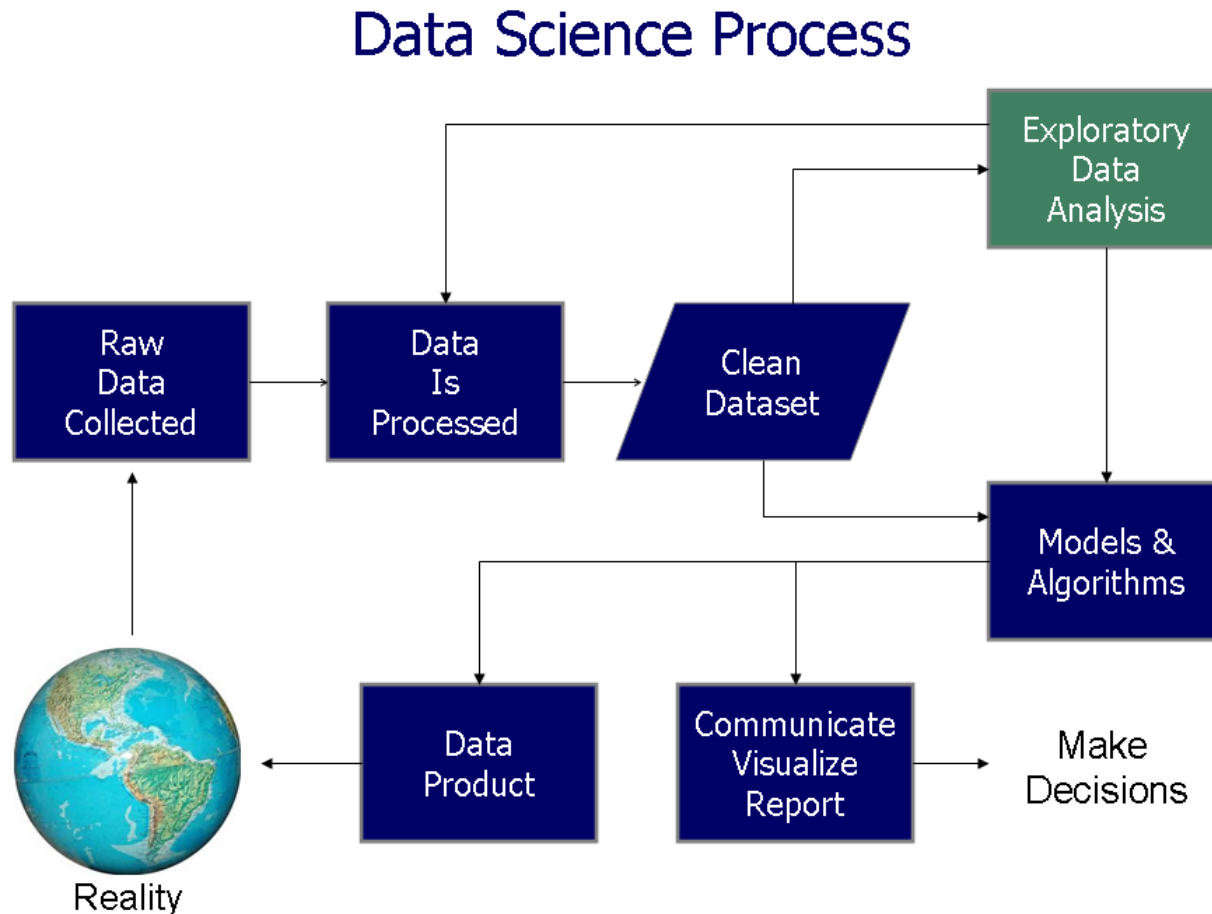  - Online Analytical Processing (OLAP)
    - Pivot Table

# What is data exploration?

- ***"A preliminary exploration of the data to better understand its characteristics"***

- Key motivations
  - Helping to select the right tool for preprocessing or analysis
  - Making use of humans' abilities to recognize patterns
    - People can recognize patterns not captured by data analysis tools

- Characteristics of data
  - Size or amount of data
  - Completeness of the data
  - Correctness of the data
  - Possible relationships amongst data elements

# Techniques Used In Data Exploration

o Summary statistics

o Visualization

o Online Analytical Processing (OLAP)

# Exploratory Data Analysis (EDA) in Data Science Process



## Data Science Process

Asst. Prof. Dr. Rachsuda Setthawong

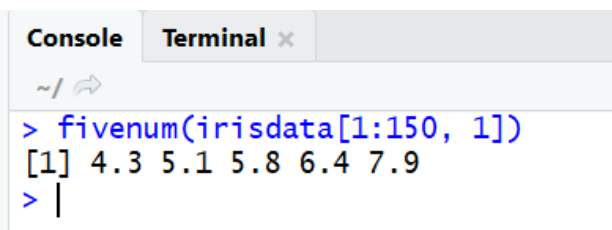https://en.wikipedia.org/wiki/Exploratory_data_analysis

# Five-number Summary in Exploratory Data Analysis (EDA)

- In statistics, John Tukey promoted the use of **five number summary** of numerical data, the two extremes (*maximum and minimum*), the median, and the quartiles in EDA.

- The **five-number summary** is a set of descriptive statistics that provide information about a dataset:

  - the sample minimum *(smallest observation)*
  - the lower quartile *(first quartile or 25th percentile)*
  - the median *(the middle value or 50th percentile)*
  - the upper quartile *( the third quartile or 75th percentile)*
  - the sample maximum *(largest observation)*

Asst. Prof. Dr. Rachsuda Setthawong

6

# Five-number Summary in R

- #Read data from file into data frame
- irisdata<-read.table("C:/Users/noox/Downloads/iris.csv", header=TRUE, sep=",")

- #use help
- help (fivenum)

- #compute 5 Tukey's 5 descriptive statistics
- fivenum(irisdata[1:150, 1]) # or fivenum(irisdata[, 1])

```
Console   Terminal ×
~/ ⇨
> fivenum(irisdata[1:150, 1])
[1] 4.3 5.1 5.8 6.4 7.9
> |
```

# Iris Sample Data Set

- Many of the exploratory data techniques are illustrated with the Iris Plant data set.
  - Source: UCI Machine Learning Repository
    https://archive.ics.uci.edu/ml/datasets/iris

  - Three flower types (classes):
    - Setosa
    - Virginica
    - Versicolour

  - Four (non-class) attributes
    - Sepal width and length
    - Petal width and length



Virginica. Robert H. Mohlenbrock. USDA NRCS. 1995. Northeast wetland flora: Field office guide to plant species. Northeast National Technical Center, Chester, PA. Courtesy of USDA NRCS Wetland Science Institute.

8

# Summary Statistics

- Summary statistics are numbers that summarize properties of the data

| Measure of **Frequency** | Measure of **Location** | Measure of **Spread** |
|---|---|---|
| Frequency | Median | Range |
| Mode | Mean | Variance |
| Percentile | | |

▢ : for categorical attributes

▢ : for continuous attributes

Asst. Prof. Dr. Rachsuda Setthawong

# Measure of **Frequency**

(Categorical attribute)

## Frequency

- The **frequency** of an attribute value is *the percentage of time the value occurs in the data set*

- *E.g.,*

## Mode

- The **mode** of an attribute is *the most frequent attribute value*

| Weight (Kg) | Frequency | Cumulative Frequency |
|---|---|---|
| 0 up to 20 | 2 | 2 |
| 20 up to 40 | 7 | 9 |
| 40 up to 60 | 12 | 21 |
| 60 up to 80 | 6 | 27 |
| 80 up to 100 | 3 | 30 |

Mode ➡ 40 up to 60

# Percentiles

(Continuous attributes)

- Given an ordinal or continuous attribute $x$ and a number $p$ between 0 and 100, the $p^{\text{th}}$ percentile is a value $X_p$ of x such that $p\%$ of the observed values of x are less than $X_p$.

- E.g., if the 80th percentile $X_{80\%}$ is 165, 80% of all values of $X$ are below *165*, and the other 20% are above *165*.

11

# Percentiles

(Continuous attributes)

The rank (R) of the 50th percentile:

**R = P/100 x (N + 1)**

R = 50/100 x (8 + 1) = 0.5 x 9 = 4.5

If R is an integer, the $P^{th}$ percentile is **the number** with rank R.

Otherwise, interpolate the $P^{th}$ percentile as follows:

1.  Define IR as the integer portion of R.

    IR = 4

2.  Define FR as the fractional portion of R.

    FR = 0.5

3.  Find the numbers (*values*) with Rank $I_R$ and with Rank $I_{R+1}$.

    $I_R = 8$

    $I_{R+1} = 9$

4.  Interpolate it:

    $$R_{Interpolate} = I_R + [\ FR * (I_{R+1} - I_R)\ ]$$

    $R_{Interpolate} = 8 + [(0.5)(9 - 8)] = 8.5$

| Number (*value*) | Rank |
|---|---|
| 3 | 1 |
| 5 | 2 |
| 7 | 3 |
| 8 | 4 |
| 9 | 5 |
| 11 | 6 |
| 13 | 7 |
| 15 | 8 |

Ref: http://onlinestatbook.com/2/introduction/percentiles.html

Asst. Prof. Dr. Rachsuda Setthawong

12

# Measures of **Location**:
## Mean and Median

- The **mean** is the most common measure of the location of a set of points.
  - Disadv: the mean is very sensitive to outliers.

$$\text{mean}(x) = \overline{x} = \frac{1}{m} \sum_{i=1}^{m} x_i$$

- The **median** (50th Percentile)

$$\text{median}(x) = \begin{cases} x_{(r+1)} & \text{if } m \text{ is odd, i.e., } m = 2r+1 \\ \frac{1}{2}(x_{(r)} + x_{(r+1)}) & \text{if } m \text{ is even, i.e., } m = 2r \end{cases}$$

- **Trimmed mean**
  - E.g., to trim the mean by 40%, we remove the lowest 20% and the highest 20% of values.

# Measures of **Spread**:
## Range and Variance

- **Range** is the difference between the max and min

- The **variance** or **standard deviation** ($\sqrt{variance}$) measures the spread of a set of points *(how far away the measurements are from the center)*.

$$\text{variance}(x) = s_x^2 = \frac{1}{m-1} \sum_{i=1}^{m} (x_i - \overline{x})^2$$

- Disadv: sensitive to outliers, so that other measures are often used.

Average Absolute Difference: $\text{AAD}(x) = \frac{1}{m} \sum_{i=1}^{m} |x_i - \overline{x}|$

Median Absolute Difference: $\text{MAD}(x) = median\left(\{|x_1 - \overline{x}|, \ldots, |x_m - \overline{x}|\}\right)$

$$\text{interquartile range}(x) = x_{75\%} - x_{25\%}$$

14

# Measures of Data Skewness

- Measures an imbalance and asymmetry from the mean of a data distribution

$$a_3 = \sum \frac{(X_i - \bar{X})^3}{ns^3}$$

s: the sample standard deviation

Note: different tools' package may use different formula

- **Normal distribution:**



Mean = Median = Mode

- **Skewness:**

•If the skewness is between -0.5 and 0.5, the data are **fairly symmetrical**

•If the skewness is between -1 and – 0.5 or between 0.5 and 1, the data are **moderately skewed**

•If the skewness is less than -1 or greater than 1, the data are **highly skewed**



Mean < Median < Mode

**Negative Skew**

(longer left-hand tail)



Mode < Median < Mean

**Positive Skew**

(longer right-hand tail)

# What does the skewness of 'Sepal Length' tell us?

| sepal length | |
|---|---:|
| | |
| Mean | 5.843333 |
| Standard Error | 0.067611 |
| Median | 5.8 |
| Mode | 5 |
| Standard Deviation | 0.828066 |
| Sample Variance | 0.685694 |
| Kurtosis | -0.55206 |
| Skewness | 0.314911 |
| Range | 3.6 |
| Minimum | 4.3 |
| Maximum | 7.9 |
| Sum | 876.5 |
| Count | 150 |
| Confidence Level(95.0%) | 0.133601 |

# Measure of **Variables' Relationship**: Covariance and Correlation

- Indicate linear relationship of 2 variables (or objects).
- Correlation also shows the degree to which the variables tend to move together.

Covariance:

$$COV(x,y) = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

Correlation: $r_{(x,y)} = \frac{COV(x,y)}{s_x s_y}$



Stock Market Returns vs. Economic Growth

Gasoline Prices vs. World Oil Production

$x$ = the independent variable
$y$ = the dependent variable
$n$ = number of data points in the sample
$\bar{x}$ = the mean of the independent variable $x$
$\bar{y}$ = the mean of the dependent variable $y$

$s_x$ = sample standard deviation of the random variable $x$
$s_y$ = sample standard deviation of the random variable $y$

$$s_x = \sqrt{\frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n-1}}$$

17

Note: for Population standard deviation, the divisor is n (*not n-1*).

Asst. Prof. Dr. Rachsuda Setthawong

# Correlation Interpretation

| Perfect Positive Correlation | High Positive Correlation | Low Positive Correlation | No Correlation | Low Negative Correlation | High Negative Correlation | Perfect Negative Correlation |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| 1 | 0.9 | 0.5 | 0 | -0.5 | -0.9 | -1 |

Correlation can have a value:

- **1** is a perfect positive correlation
- **0** is no correlation (the values don't seem linked at all)
- **-1** is a perfect negative correlation

Asst. Prof. Dr. Rachsuda Setthawong

# Tools for Summary Statistics

- R
- Python
- Minitab
- SPSS
- MS Excel
- Weka
- Rapidminer
- Etc.

Tableau:
https://www.tableau.com/trial/data-visualization?utm_campaign_id=2017049&utm_campaign=Prospecting-CORE-ALL-ALL-ALL-ALL&utm_medium=Paid+Search&utm_source=Google+Search&utm_language=EN&utm_country=SEA&kw=%2Bvisualization%20%2Btableau&adgroup=CTX-Brand-Data+Visualization-EN-B&adused=324827239180&matchtype=b&placement=&gclid=EAIaIQobChMIocCWpKuS6gIVyhErCh0G6QaIEAAYASAAEgJ5CPD_BwE&gclsrc=aw.ds

# Outlines

- Data Exploration: Tasks and Techniques
  - Summary Statistics
  - Data Visualization
  - Online Analytical Processing (OLAP)
    - Pivot Table

# Visualization

- **The conversion of data into a visual or tabular format** so that the characteristics of the data and the relationships among data items or attributes can be analyzed or reported.

- One of the most powerful and appealing techniques for data exploration:
  - Humans have a well-developed ability to analyze large amounts of information that is presented visually
  - Can detect general patterns and trends
  - Can detect outliers and unusual patterns

# Example: Sea Surface Temperature

- The following shows the Sea Surface Temperature (SST) for July 1982
  - Tens of thousands of data points are summarized in a single figure

22

# Representation

- Data (objects, their attributes, and their relationships) are translated into graphical elements (points, lines, shapes, and colors).

# Arrangement

- Is the placement of visual elements within a display

- Can make a large difference in how easy it is to understand the data

- Example:

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 1 | 1 | 0 |
| 2 | 1 | 0 | 1 | 0 | 0 | 1 |
| 3 | 0 | 1 | 0 | 1 | 1 | 0 |
| 4 | 1 | 0 | 1 | 0 | 0 | 1 |
| 5 | 0 | 1 | 0 | 1 | 1 | 0 |
| 6 | 1 | 0 | 1 | 0 | 0 | 1 |
| 7 | 0 | 1 | 0 | 1 | 1 | 0 |
| 8 | 1 | 0 | 1 | 0 | 0 | 1 |
| 9 | 0 | 1 | 0 | 1 | 1 | 0 |

|   | 6 | 1 | 3 | 2 | 5 | 4 |
|---|---|---|---|---|---|---|
| 4 | 1 | 1 | 1 | 0 | 0 | 0 |
| 2 | 1 | 1 | 1 | 0 | 0 | 0 |
| 6 | 1 | 1 | 1 | 0 | 0 | 0 |
| 8 | 1 | 1 | 1 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 1 | 1 | 1 |
| 3 | 0 | 0 | 0 | 1 | 1 | 1 |
| 9 | 0 | 0 | 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 | 1 | 1 | 1 |
| 7 | 0 | 0 | 0 | 1 | 1 | 1 |

24

# Selection

- Choosing a subset of attributes
  - Dimensionality reduction is often used to reduce the number of dimensions to two or three

- Choosing a subset of objects
  - A region of the screen can only show limited no. of points
  - Can sample, but want to preserve points in sparse areas

# Visualization Techniques: Histograms

- **Histogram**
  - Usually **shows the distribution of values of a single variable**
  - Divide the values into **bins** and show a bar plot of the number of objects in each bin.
  - The **height** of each bar indicates the number of objects
  - Shape of histogram depends on the number of bins
- Example: Petal Width (10 and 20 bins, respectively)

# Two-Dimensional Histograms

- Show the joint distribution of the values of two attributes

- Example: petal width and petal length
  - What does this tell us?

# Visualization Techniques: Box Plots

- Another way of displaying the distribution of data



- outlier
- 90th percentile
- 75th percentile
- 50th percentile (Median)
- 25th percentile
- 10th percentile

How to create Box Plot: https://www.youtube.com/watch?v=ucWmfmXb1kk

28

# Example of Box Plots

- Box plots can be used to compare attributes

https://medium.com/dayem-siddiqui/understanding-and-interpreting-box-plots-d07aab9d1b6c

# Visualization Techniques: Scatter Plots

- Attributes values determine the position

- Two-dimensional scatter plots most common, but can have three-dimensional scatter plots

- Often additional attributes can be displayed by using the size, shape, and color of the markers that represent the objects

- It is useful to have arrays of scatter plots can compactly summarize the relationships of several pairs of attributes
  - See example on the next slide

# Scatter Plot Array of Iris Attributes

31

# Visualization Techniques: Contour Plots

- Useful when a continuous attribute is measured on a spatial grid.

- Partition the plane into regions of similar values.

- The contour lines that form the boundaries of these regions connect points with equal values.

32

# Contour Plot Example:

Sea Surface Temperature (Dec, 1998)



Celsius

# Visualization Techniques: Matrix Plots

- Can plot the data matrix

- Useful when objects are sorted according to class

- Typically, the attributes are normalized to prevent one attribute from dominating the plot

# Visualization of the Iris Correlation Matrix

35

# Visualization Techniques: Parallel Coordinates

- Used to plot the attribute values of high-dimensional data

- Use a set of parallel axes instead of using perpendicular axes.

- Thus, **each object is represented as a line**
  - The attribute values of each object are plotted as a point on each corresponding coordinate axis, and
  - The points are connected by a line

36

# Parallel Coordinates Plots for Iris Data



Can rearrange order of attributes

# Other Visualization Techniques

- Star Plots
  - Similar approach to parallel coordinates, but axes radiate from a central point
  - The line connecting the values of an object is a polygon

- Chernoff Faces
  - Approach created by Herman Chernoff
  - This approach associates each attribute with a characteristic of a face
  - The values of each attribute determine the appearance of the corresponding facial characteristic
  - Each object becomes a separate face
  - Relies on human's ability to distinguish faces

# Star Plots for Iris Data

Setosa

1  2  3  4  5

Versicolour

51  52  53  54  55

Virginica

101  102  103  104  105

# Chernoff Faces for Iris Data

Setosa

Versicolour

Virginica

40

# Sample Charts in MS Excel

41

Microsoft Tutorial – Creating a Graph from Start to Finish:
https://support.office.com/en-us/article/create-a-chart-from-start-to-finish-0baf399e-dd61-4e18-8a73-b3fd5d5680c2
Juicebox Chart Chooser: http://labs.juiceanalytics.com/chartchooser/

# Outlines

- Data Exploration: Tasks and Techniques
  - Summary Statistics
  - Data Visualization
  - Online Analytical Processing (OLAP)
    - Pivot Table

Asst. Prof. Dr. Rachsuda Setthawong

# On-Line Analytical Processing (OLAP)

- Relational databases put data into tables, while OLAP uses a multidimensional array representation.
  - Such representations of data previously existed in statistics and other fields

- There are a number of data analysis and data exploration operations that are easier with such a data representation.

# Creating a Multidimensional Array

- Two steps (Tabular data to a multidimensional array):

  - First, identify which attributes are to be the dimensions and which attribute is to be the target attribute whose values appear as entries in the multidimensional array.

    - The attributes used as **dimensions** must have **discrete values**

    - The **target value** is typically a **count** or **continuous value**, e.g., the cost of an item

    - Can have no target variable at all except the count of objects that have the same set of attribute values

  - Second, find the value of each entry in the multidimensional array by summing the values (of the target attribute) or count of all objects that have the attribute values corresponding to that entry.

# Example: Iris data

1.  **Define dimensions:** A) petal width and B) petal length, and C) Species Type

    *   **Discretize** A and B to have categorical values: *low*, *medium*, and *high*

2.  **Define target attribute:** **count** attribute and count values of each case

| Petal Length | Petal Width | Species Type | Count |
|:---:|:---:|:---:|:---:|
| low | low | Setosa | 46 |
| low | medium | Setosa | 2 |
| medium | low | Setosa | 2 |
| medium | medium | Versicolour | 43 |
| medium | high | Versicolour | 3 |
| medium | high | Virginica | 3 |
| high | medium | Versicolour | 2 |
| high | medium | Virginica | 3 |
| high | high | Versicolour | 2 |
| high | high | Virginica | 44 |

# Example: Iris data (continued)

- Each unique tuple of petal width, petal length, and species type identifies <span style="color:red">one element</span> of the array.

- This element is assigned the corresponding <span style="color:blue">count</span> value.

- All non-specified tuples are 0.

| Petal Length | Petal Width | Species Type | Count |
|---|---|---|---|
| low | low | Setosa | 46 |
| low | medium | Setosa | 2 |
| medium | low | Setosa | 2 |
| medium | medium | Versicolour | 43 |
| medium | high | Versicolour | 3 |
| medium | high | Virginica | 3 |
| high | medium | Versicolour | 2 |
| high | medium | Virginica | 3 |
| high | high | Versicolour | 2 |
| high | high | Virginica | 44 |

46

# Example: Iris data (continued)

- Slices of the multidimensional array are shown by the following cross-tabulations

| Petal Length | Petal Width | Species Type | Count |
|---|---|---|---|
| low | low | Setosa | 46 |
| low | medium | Setosa | 2 |
| medium | low | Setosa | 2 |
| medium | medium | Versicolour | 43 |
| medium | high | Versicolour | 3 |
| medium | high | Virginica | 3 |
| high | medium | Versicolour | 2 |
| high | medium | Virginica | 3 |
| high | high | Versicolour | 2 |
| high | high | Virginica | 44 |

### Setosa

| Length \ Width | low | medium | high |
|---|---|---|---|
| low | 46 | 2 | 0 |
| medium | 2 | 0 | 0 |
| high | 0 | 0 | 0 |

### Versicolour

| Length \ Width | low | medium | high |
|---|---|---|---|
| low | 0 | 0 | 0 |
| medium | 0 | 43 | 3 |
| high | 0 | 2 | 2 |

### Virginica

| Length \ Width | low | medium | high |
|---|---|---|---|
| low | 0 | 0 | 0 |
| medium | 0 | 0 | 3 |
| high | 0 | 3 | 44 |

# OLAP Operations: Data Cube

- The key operation of a OLAP is the formation of a data cube

- A data cube is **a multidimensional representation of data**, together with all possible aggregates.

- By all possible aggregates, we mean the aggregates that result by **selecting a proper subset of the dimensions** and **summing over all remaining dimensions.**

- For example, if we choose the species type dimension of the Iris data and sum over all other dimensions, the result will be a one-dimensional entry with three entries, each of which gives the number of flowers of each type.

# Data Cube Example

- Consider a data set that records the sales of products at a number of company stores at various dates.

- This data can be represented as a three-dimensional array

- There are 3 two-dimensional aggregates (3 choose 2), 3 one-dimensional aggregates, and 1 zero-dimensional aggregate (the overall total)



49

# Data Cube Example (continued)

- The following figure table shows one of the two dimensional aggregates, along with two of the one-dimensional aggregates, and the overall total

50

# OLAP Operations:
## Slicing and Dicing

- Slicing is **selecting a group of cells** from the entire multidimensional array **by specifying a specific value for one or more dimensions.**

E.g., Slice 1 quarter

51

# OLAP Operations:
## Slicing and Dicing

- **Dicing** involves **selecting a subset of cells** by specifying a **range** **of attribute values**.
  - This is equivalent to defining a subarray from the complete array.



- In practice, both operations can also be accompanied by aggregation over some dimensions.
  https://cracklogic.com/olap-tutorial/#Slice

# OLAP Operations: Roll-up and Drill-down

- **Attribute values** often have a **hierarchical structure**.
  - Each *date* is associated with a year, month, and week.
  - A *location* is associated with a continent, country, state (province, etc.), and city.
  - *Products* can be divided into various categories, such as clothing, electronics, and furniture.

- Note that these categories often nest and form a tree
  - A year contains months which contains day
  - A country contains a state which contains a city

# OLAP Operations: Roll-up and Drill-down

- This **hierarchical structure** gives rise to the roll-up and drill-down operations.

# Outlines

- Data Exploration: Tasks and Techniques
  - Summary Statistics
  - Data Visualization
  - Online Analytical Processing (OLAP)
    - Pivot Table

Asst. Prof. Dr. Rachsuda Setthawong

# Pivot Table in MS Excel

56

https://exceljet.net/excel-pivot-tables

# What is Pivot Table?

- A **table** of statistics that summarizes the data of a more extensive **table** (e.g., a database, spreadsheet).

- This summary include sums, averages, or other statistics, which the **pivot table** groups together in a meaningful way.

### Raw Data

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | Order ID | Product | Category | Amount | Date | Country | | |
| 2 | 1 | Carrots | Vegetables | $4,270 | 1/6/2016 | United States | | |
| 3 | 2 | Broccoli | Vegetables | $8,239 | 1/7/2016 | United Kingdom | | |
| 4 | 3 | Banana | Fruit | $617 | 1/8/2016 | United States | | |
| 5 | 4 | Banana | Fruit | $8,384 | 1/10/2016 | Canada | | |
| 6 | 5 | Beans | Vegetables | $2,626 | 1/10/2016 | Germany | | |
| 7 | 6 | Orange | Fruit | $3,610 | 1/11/2016 | United States | | |
| 8 | 7 | Broccoli | Vegetables | $9,062 | 1/11/2016 | Australia | | |
| 9 | 8 | Banana | Fruit | $6,906 | 1/16/2016 | New Zealand | | |
| 10 | 9 | Apple | Fruit | $2,417 | 1/16/2016 | France | | |
| 11 | 10 | Apple | Fruit | $7,431 | 1/16/2016 | Canada | | |

### Pivot Tables

| | A | B | C |
|---|---|---|---|
| 1 | Country | (All) | |
| 2 | | | |
| 3 | Row Labels | Sum of Amount | |
| 4 | Banana | 340295 | |
| 5 | Apple | 191257 | |
| 6 | Broccoli | 142439 | |
| 7 | Carrots | 136945 | |
| 8 | Orange | 104438 | |
| 9 | Beans | 57281 | |
| 10 | Mango | 57079 | |
| 11 | Grand Total | 1029734 | |
| 12 | | | |

| | A | B | C |
|---|---|---|---|
| 1 | Country | France | |
| 2 | | | |
| 3 | Row Labels | Count of Amount | |
| 4 | Apple | 16 | |
| 5 | Banana | 7 | |
| 6 | Carrots | 1 | |
| 7 | Mango | 1 | |
| 8 | Orange | 1 | |
| 9 | Beans | 1 | |
| 10 | Broccoli | 1 | |
| 11 | Grand Total | 28 | |
| 12 | | | |

| | A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Category | (All) | | | | | | | | |
| 2 | | | | | | | | | | |
| 3 | Sum of Amount | Column | | | | | | | | |
| 4 | Row Labels | Apple | Banana | Beans | Broccoli | Carrots | Mango | Orange | Grand Total | |
| 5 | Australia | 20634 | 52721 | 14433 | 17953 | 8106 | 9186 | 8680 | 131713 | |
| 6 | Canada | 24867 | 33775 | | 12407 | | 3767 | 19929 | 94745 | |
| 7 | France | 80193 | 36094 | 680 | 5341 | 9104 | 7388 | 2256 | 141056 | |
| 8 | Germany | 9082 | 39686 | 29905 | 37197 | 21636 | 8775 | 8887 | 155168 | |
| 9 | New Zealand | 10332 | 40050 | | 4390 | | | 12010 | 66782 | |
| 10 | United Kingdom | 17534 | 42908 | 5100 | 38436 | 41815 | 5600 | 21744 | 173137 | |
| 11 | United States | 28615 | 95061 | 7163 | 26715 | 56284 | 22363 | 30932 | 267133 | |
| 12 | Grand Total | 191257 | 340295 | 57281 | 142439 | 136945 | 57079 | 104438 | 1029734 | |
| 13 | | | | | | | | | | |

Asst. Prof. Dr. Rachsuda Setthawong

57

# Creating Pivot Table: 1/6

1. Open the file **SalesData.xlsx**
2. Add PivotTable

# Creating Pivot Table: 2/6
## Selecting a Table or Range of Data

Asst. Prof. Dr. Rachsuda Setthawong

## The Worksheet with Pivot Table Added

# Creating Pivot Table: 4/6
## Selecting Fields to be Displayed in Pivot Table

# Creating Pivot Table: 5/6
## Changing Summary Info.

Asst. Prof. Dr. Rachsuda Setthawong

# Creating Pivot Table: 6/6
## Rearranging Fields in Rows and Changing Summary Info.



Can be expanded for detailed data

# Example of Summary 1:
## Order Amount per Salesperson in Different Countries

Asst. Prof. Dr. Rachsuda Setthawong

# Example of Summary 2 (using Filter):
## Order Amount in 2003 per Salesperson in Different Countries

# Example of Summary 3 (Aggregation using Group):
Grouping Order Amount by Year per Salesperson in Different Countries

# Example of Summary 3 (Aggregation using Group): 1/3

# Example of Summary 3 (Aggregation using Group): 3/3

# Example 4: Slicing – 1/3

# Example 4: Slicing – 2/3

# Example 4: Slicing – 3/3



Select subset of 'Order Date' will update the Pivot Table to show only the selected month

# Creating Pivot Chart

Asst. Prof. Dr. Rachsuda Setthawong

# Pivot Chart Added

Asst. Prof. Dr. Rachsuda Setthawong

# Adding Fields to the Chart

# More Action on Chart:
## Interactively Update the Chart using Filter

# More Action on Chart:
## Interactively Update the Chart using Filter: Result

# More Action on Pivot Table (Chart):
## Using Timeline to Filter Data
### (Must be 'Date/Time' Field)



[2]

[3]

[4]

Insert Timelines

Order Date

[1]

[6]
Data in the pivot Table (and graph if any) are updated, accordingly.

[5]
Select JAN – JULY

# More Action on Pivot Table (Chart):
## Changing Timeline's Granularity

# Analysis ToolPak in MS Excel

Asst. Prof. Dr. Rachsuda Setthawong

# Loading and the Activating Analysis ToolPak

1. Click the **File** tab, click **Options**, and then click the **Add-Ins** category.

2. If you're using Excel 2007, click the **Microsoft Office Button** , and then click **Excel Options**

3. In the **Manage** box, select **Excel Add-ins** and then click **Go**.

4. If you're using Excel for Mac, in the file menu go to **Tools** > **Excel Add-ins**.

5. In the **Add-Ins** box, check the **Analysis ToolPak** check box, and then click **OK**.
   - If **Analysis ToolPak** is not listed in the **Add-Ins available** box, click **Browse** to locate it.

   - If you are prompted that the Analysis ToolPak is not currently installed on your computer, click **Yes** to install it.

# Click the **Data** tab,
## in Analysis section click **Data Analysis**

- Anova
- Correlation
- Covarience
- Descriptive Statistics
- Exponential Smoothing
- F-Test Two-Sample for Variances
- Fourier Analysis
- Histogram
- Moving Average
- Random Number Generation
- Rank and Percentile
- Regression
- Sampling
- t-Test
- z-Test