# BUSINESS INTELLIGENCE

LECTURE 1

# OUTLINE

- Introduction
- Data mining overview
- Data pre-processing
- Classification techniques
- Clustering techniques
- Association rules
- References

# INTRODUCTION

# ALICE'S STORY

- Here is Alice

- She is interested in F&B

- One day, she opened her first store in Saigon

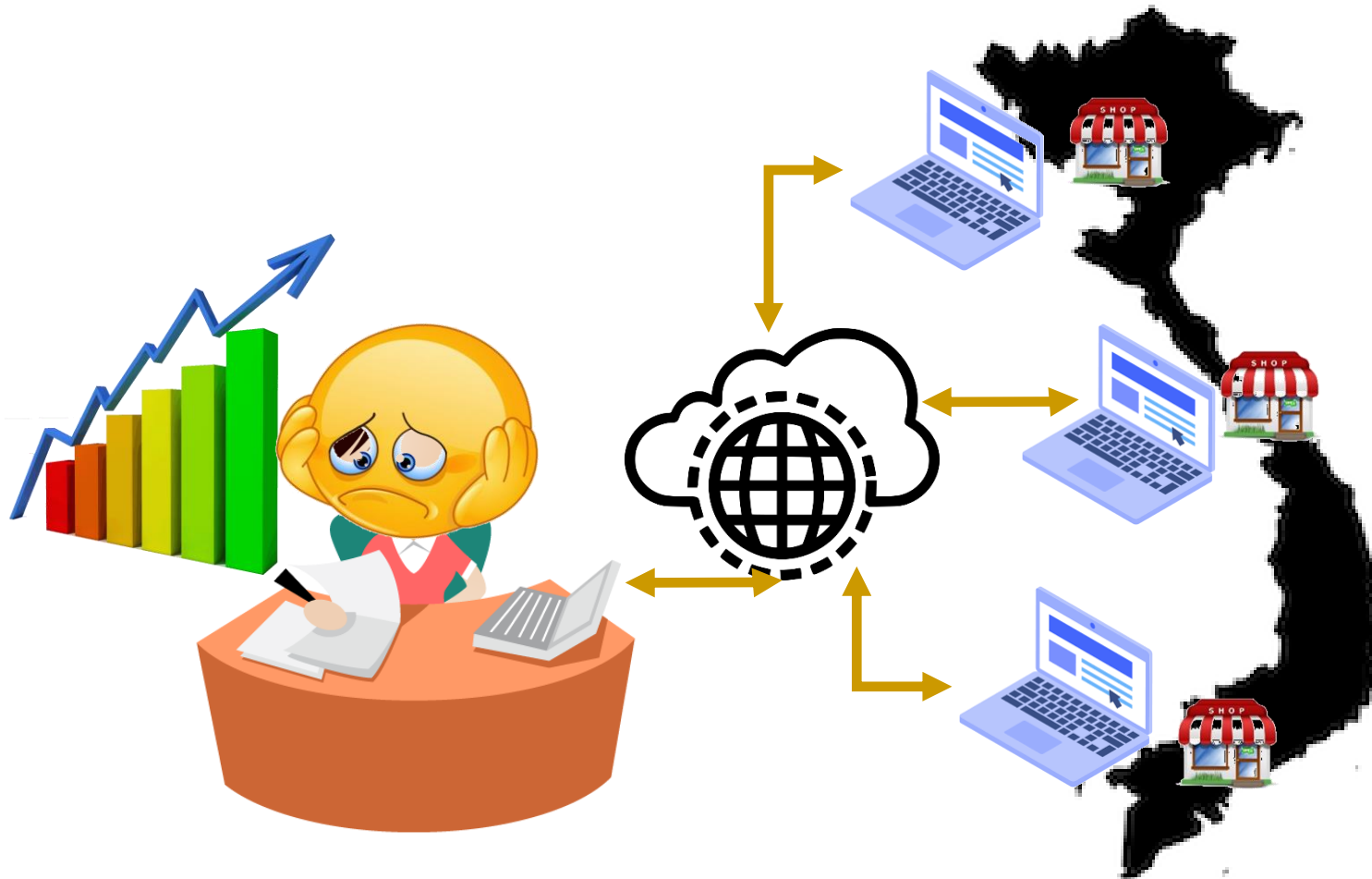- It's a traditional restaurant with Vietnamese cuisine
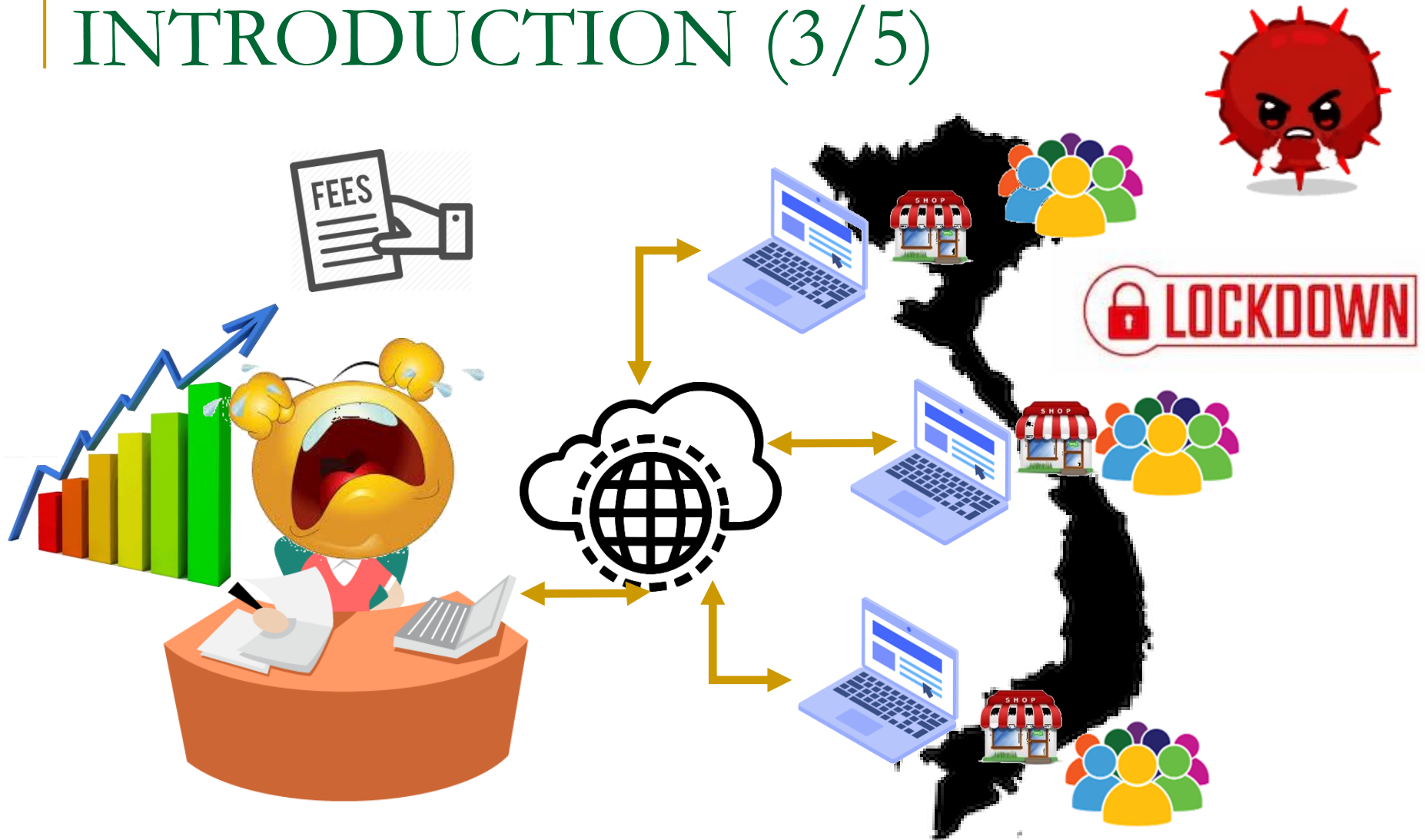
# INTRODUCTION (1/5)

# INTRODUCTION (3/5)

# DISCUSSION

# WOULD YOU SUGGEST ANY SOLUTIONS TO HELP HER?

# INTRODUCTION (4/5)

# INTRODUCTION (5/5)

# AND ANOTHER PROBLEM COMES...

- It takes her long time to summarize sales measurements
- She has to redo all for the new cycle (i.e., week, month, and year)
- When there is any error, she has to update her data and recheck the numbers
- And more...

SUCH A RICE CONSUMPTION SYSTEM!

# DISCUSSION

# SHE NEEDS YOUR SUPPORT AGAIN!

# YOU MAY COME UP WITH...



Query

# OR BETTER FOR ALICE...

15

# DASHBOARD

# SALES FORECAST

# ANOMALY DETECTION



https://miro.medium.com/max/1400/1*7qWjee_M6PTrAd0PJdU1yg.png

# CLUSTER IDENTIFICATION



https://miro.medium.com/max/1400/1*sUeqYcENVII1RlyYA_-Uxg.png

# BI SYSTEM

- "Business intelligence is an infrastructure that helps in the process of collecting, storing, and analyzing data from business operations.

- BI provides comprehensive business metrics, in near-real-time, to support better decision making. You can create performance benchmarks, spot market trends, increase compliance, and improve almost every aspect of your business with better business intelligence."

https://www.tableau.com/learn/articles/business-intelligence/bi-business-analytics

# BI CHALLENGES

- From business
  - Reports meet the needs of business requirements and different stakeholders
  - Report usability
  - More than expectation
  - Etc.
- From technique and technology
  - BI system development
  - Different levels of BI
  - The support of BI technology
  - Etc.

# WHAT'S MORE…

- **Lack of credibility**
  - Sales says revenue is +10%
  - Finance says revenue is -15%
  - SCM says revenue is +5%
- **Reasons**
  - No time basis of data
  - The algorithmic differential of data
  - The levels of extraction
  - The problem of external data
  - No common source of data from the beginning
  - Errors of all kinds
  - Etc.

[1]



Department A +10%

Department B -15%

22

# WHAT'S MORE…

- ## Problems of productivity
  - Locate and analyze the data for the report.
  - Compile the data for the report.
  - Get programmer/analyst resources to accomplish these two tasks.

- ## Reasons
  - Big data (3Vs)
  - Technology variance
  - Reusability
  - Unknown future demands
  - Etc.

[1]



PRODUCTIVITY

Produce a corporate report, across all data.

Locating the data requires looking at lots of files.

Lots of extract programs, each customized, have to cross many technological barriers.

# WHAT'S MORE…

- **Problems of obtaining information**
  - Data is unavailable
  - Data is not integrated
- **Reasons**
  - No data history
  - Unforeseen data needs
  - Unintegrated legacy applications
  - Etc.



GOING FROM DATA TO INFORMATION

Loans
DDA
CD
Passbook

First, you run into lots of applications.

Loans
DDA
CD
Passbook
Same element, different name
Different element, same name
Element exists only here

Next, you run into the lack of integration across applications.

[1]

# ARCHITECTED ENVIRONMENT

**LEVELS OF THE ARCHITECTURE**



Operational → Atomic/data warehouse → Departmental → Individual

| Operational | Atomic/data warehouse | Departmental | Individual |
|---|---|---|---|
| • Detailed | • Most granular | • Parochial | • Temporary |
| • Day to day | • Time variant | • Some derived; | • Ad hoc |
| • Current valued | • integrated | some primitive | • Heuristic |
| • High probability of access | • Subject-oriented | • Typical departments | • Non-repetitive |
| • Application-oriented | • some summary | • Accounts | • PC, work-station based |
| | | • Marketing | |
| | | • Engineering | |
| | | • Actuarial | |
| | | • Manufacturing | |

**Figure 1-10**  Although it is not apparent at first glance, there is very little redundancy data across the architected environment.

## A CHANGE IN APPROACHES

| PRIMITIVE DATA/OPERATIONAL DATA | DERIVED DATA/DSS DATA |
|---|---|
| • Application-oriented | • Subject-oriented |
| • Detailed | • Summarized, otherwise refined |
| • Accurate, as of the moment of access | • Represents values over time, snapshots |
| • Serves the clerical community | • Serves the managerial community |
| • Can be updated | • Is not updated |
| • Run repetitively | • Run heuristically |
| • Requirements for processing understood *a priori* | • Requirements for processing not understood *a priori* |
| • Compatible with the SDLC | • Completely different life cycle |
| • Performance-sensitive | • Performance relaxed |
| • Accessed a unit at a time | • Accessed a set at a time |
| • Transaction-driven | • Analysis-driven |
| • Control of update a major concern in terms of ownership | • Control of update no issue |
| • High availability | • Relaxed availability |
| • Managed in its entirety | • Managed by subsets |
| • Nonredundancy | • Redundancy is a fact of life |
| • Static structure; variable contents | • Flexible structure |
| • Small amount of data used in a process | • Large amount of data used in a process |
| • Supports day-to-day operations | • Supports managerial needs |
| • High probability of access | • Low, modest probability of access |

**Figure 1-9**  Differences between primitive and derived data.

[1]

# AN EXAMPLE

**A SIMPLE EXAMPLE–A CUSTOMER**



| Operational | Atomic/data warehouse | Dept/data mart customers by month | Individual |
|---|---|---|---|
| J Jones<br>123 Main St<br>Credit - AA | J Jones<br>1986-1987<br>456 High St<br>Credit - B | Jan - 4101<br>Feb - 4209<br>Mar - 4175<br>Apr - 4215<br>. . . . . . . . . . . .<br>. . . . .<br>. . . . . . . . . . . .<br>. . . . . | customers since 1982 with acct balances > 5,000 and with credit ratings of B or higher |
| | J Jones<br>1987-1989<br>456 High St<br>Credit - A | | Temporary! |
| | J Jones<br>1989-pres<br>123 Main St<br>Credit - AA | | |
| What is J Jones's credit rating right now? | What has been the credit history of J Jones? | Are we attracting more or fewer customers over time? | What trends are there for the customers we are analyzing? |

**Figure 1-11** The kinds of queries for which the different levels of data can be used.

[1]

26

# TRADITIONAL TO MODERN BI (1/4)



Deloitte: Modern Business Intelligence, The Path to Big Data Analytics, 2018.

# TRADITIONAL TO MODERN BI (2/4)

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| **1** | Name | HP | Attack | Defense | Sp. Atk | Sp. Def | Speed | |
| **2** | Bulbasaur | 45 | 49 | 49 | 65 | 65 | 45 | |
| **3** | Charmander | 39 | 52 | 43 | 60 | 50 | 65 | |
| **4** | Squirtle | 44 | 48 | 65 | 50 | 64 | 43 | |

**Chart Title**

(Bar chart comparing Bulbasaur, Charmander, and Squirtle across HP, Attack, Defense, Sp. Atk, Sp. Def, Speed. Vertical axis from 0 to 70.)

Legend: ■ Bulbasaur ■ Charmander ■ Squirtle

https://www.w3schools.com/excel/img_excel_chart_intro_8_5.png

28

# TRADITIONAL TO MODERN BI (3/4)



https://www.tableau.com/learn/articles/business-intelligence

https://www.holistics.io/

# TRADITIONAL TO MODERN BI (4/4)



https://www.tableau.com/learn/articles/business-intelligence

# BI PROCESSES AND ACTIVITIES

- Reporting and visualization
- Statistical analysis
- Data mining
- Performance metrics and benchmarking
- Visual storytelling
- Data preparation
- Querying
- Etc.

# BI STRATEGY BY TABLEAU

1. Know your business strategy and goals.
2. Identify key stakeholders.
3. Choose a sponsor from your key stakeholders.
4. Choose your BI platform and tools.
5. Create a BI team.
6. Define your scope.
7. Prepare your data infrastructure.
8. Define your goals and roadmap.

https://www.tableau.com/learn/articles/business-intelligence

# DISCUSSION

# BUSINESS ANALYTICS
# VS.
# DATA ANALYTICS

# DISCUSSION

## BI PROS AND CONS

# DISCUSSION

# BI APPLICATION
# (E.g., PowerBI, Tableau, Looker Studio, Holistics)

# DATA MINING OVERVIEW

# KNOWLEDGE DISCOVERY

**Figure 1.4** Data mining as a step in the process of knowledge discovery.

38

# NORMALIZATION VS. DE-NORMALIZATION

- **Normalization** is used to decompose one table data into different sub-tables data

- **Denormalization** is used to combine multiple table data into one

# OLTP vs. OLAP

**Table 4.1**  Comparison of OLTP and OLAP Systems

| Feature | OLTP | OLAP |
|---|---|---|
| Characteristic | operational processing | informational processing |
| Orientation | transaction | analysis |
| User | clerk, DBA, database professional | knowledge worker (e.g., manager, executive, analyst) |
| Function | day-to-day operations | long-term informational requirements decision support |
| DB design | ER-based, application-oriented | star/snowflake, subject-oriented |
| Data | current, guaranteed up-to-date | historic, accuracy maintained over time |
| Summarization | primitive, highly detailed | summarized, consolidated |
| View | detailed, flat relational | summarized, multidimensional |
| Unit of work | short, simple transaction | complex query |
| Access | read/write | mostly read |
| Focus | data in | information out |
| Operations | index/hash on primary key | lots of scans |
| Number of records accessed | tens | millions |
| Number of users | thousands | hundreds |
| DB size | GB to high-order GB | $\geq$ TB |
| Priority | high performance, high availability | high flexibility, end-user autonomy |
| Metric | transaction throughput | query throughput, response time |

*Note:* Table is partially based on Chaudhuri and Dayal [CD97].

[8]

# DATA MINING

- Data mining is the process of discovering interesting patterns and knowledge from large amounts of data.



**Figure 1.11** Data mining adopts techniques from many domains.

# DATA MINING FUNCTIONALITY

- ## Descriptive tasks
  - characterize properties of the data in a target data set

- ## Predictive tasks
  - perform induction on the current data in order to make predictions

# KINDS OF DATA TO BE MINED

- Database data

- Data warehouse data

- Transactional data

- Stream data

- Sequence data

- Graph data

- Spatial data

- Text data

- Multimedia data

- Etc.

# KINDS OF PATTERNS TO BE MINED

- Class description
- Frequent patterns
- Associations
- Classification and regression
- Cluster analysis
- Outlier analysis

# DATA CHARACTERIZATION

- It summarizes the data of the class under study (often called the target class) in general terms

- E.g.,

  - Summarizing the characteristics of customers who spend more than $5000 a year at AllElectronics.

  - The result is a general profile of these customers, such as that they are 40 to 50 years old, employed, and have excellent credit ratings.

  - The data mining system should allow the customer relationship manager to drill down on any dimension, such as on occupation to view these customers according to their type of employment.

# DATA DISCRIMINATION

- It compares the target class with one or a set of comparative classes (often called the contrasting classes)
- E.g.,
  - Comparing two groups of customers—those who shop for computer products regularly (e.g., more than twice a month) and those who rarely shop for such products (e.g., less than three times a year).
  - The resulting description provides a general comparative profile of these customers, such as that 80% of the customers who frequently purchase computer products are between 20 and 40 years old and have a university education, whereas 60% of the customers who infrequently buy such products are either seniors or youths, and have no university degree. Drilling down on a dimension like occupation, or adding a new dimension like income level, may help to find even more discriminative features between the two classes.

# FREQUENT PATTERNS

- They are patterns that occur frequently in data
  - Frequent itemsets: a set of items that often appear together in a transactional data set
  - Frequent subsequences: a frequently occurring subsequence
  - Frequent substructures: a frequently occurring structural forms (e.g., graphs, trees, or lattices) that may be combined with itemsets or subsequences.

# ASSOCIATIONS

- Mining frequent patterns leads to the discovery of interesting associations within data.

- E.g.,
  - *buys(X, "computer")* → *buys(X, "software")* with *support = 1%* and *confidence = 50%*
  - *age(X, "20..29")* and *income(X, "40K..49K")* → *buys(X, "laptop")* *with support = 2%* and *confidence = 60%*

# CLASSIFICATION

■ Classification is the process of finding a model (or function) that describes and distinguishes data classes or concepts.

age(X, "youth") AND income(X, "high")  ⟶  class(X, "A")
age(X, "youth") AND income(X, "low")  ⟶  class(X, "B")
age(X, "middle_aged")  ⟶  class(X, "C")
age(X, "senior")  ⟶  class(X, "C")

**(a)**



**Figure 1.9** A classification model can be represented in various forms: (a) IF-THEN rules, (b) a decision tree, or (c) a neural network.

[2]

# REGRESSION

- Regression models continuous-valued functions. That is, regression is used to predict missing or unavailable numerical data values rather than (discrete) class labels.

Simple linear regression:

$$Y = a + bX + u$$

Multiple linear regression:

$$Y = a + b_1 X_1 + b_2 X_2 + b_3 X_3 + ... + b_t X_t + u$$

**where:**

$Y = $ The dependent variable you are trying to predict or explain

$X = $ The explanatory (independent) variable(s) you are using to predict or associate with Y

$a = $ The y-intercept

$b = $ (beta coefficient) is the slope of the explanatory variable(s)

$u = $ The regression residual or error term

Internet

# CLUSTER ANALYSIS

- Clusters of objects are formed so that objects within a cluster have high similarity in comparison to one another, but are rather dissimilar to objects in other clusters.

- Clustering analyzes data objects without consulting class labels.

- E.g., interests based on locations

# OUTLIER ANALYSIS

- Objects that do not comply with the general behavior or model of the data. These data objects are outliers.

- E.g., fraudulent usage of credit cards

# DISCUSSION



# ARE ALL PATTERNS INTERESTING?

# DISCUSSION

# DATA MINING VS. MACHINE LEARNING?

# TO BE CONTINUED…

- Lecture 2: Data pre-processing

# QUESTIONS AND ANSWERS



Picture from: http://philadelphiasculpturegym.blogspot.com/2013/09/save-date-free-talk-and-q-on-affordable.html

# REFERENCES

[1] Tobias Zwingmann, AI-Powered Business Intelligence, Kindle Edition, O'reilly Press, 2022

[2] Jiawei Han, Micheline Kamber, "Data Mining: Concepts and Techniques", Third Edition, Morgan Kaufmann Publishers, 2012.

[3] Jeen Su Lim, John Heinrichs, "Digital Business Intelligence Management with Big Data Analytics" Kindle Edition, O'reilly Press, 2021.

[4] David L. Olson, Dursun Delen, "Advanced Data Mining Techniques", Springer-Verlag, 2008.

[5] Brian Larson, "Delivering Business Intelligence with Microsoft SQL Server 2016", McGraw-Hill Education; 4 edition, 2016.

[6] Oracle, "Data Mining Concepts", 18c, E83730-03, 2018

[7] Oracle, "Data Mining Application Developer's Guide", 2013.