

Data Warehouse & Data Mining

Chapter 1

Introduction



• Why Data Mining?

Data Mining = การทำเหมืองข้อมูล / ขุดเหมืองข้อมูล

ในปัจจุบัน Data เข้มข้น เร็วขึ้น และสามารถเก็บได้

เช่น เซ็นเซอร์ อุตสาหกรรม, เซ็นเซอร์ ปริมาณน้ำฝน \Rightarrow ซิงค์เครื่องมือที่มี การเก็บข้อมูลไปเรื่อยๆ

ดังนั้น ข้อมูลก็จะมีการ เพิ่มขึ้นเรื่อยๆ ทำให้ ล้น !!

รวมไปถึง website เช่น สหภาพทำ กับสินค้า ที่ไปดู Website เจ้าของเว็บไซต์ สามารถที่จะดู Website นั้นได้ สามารถ ทำนายได้ ว่ามีคนไปดูเยอะ ในส่วนนั้นของเว็บไซต์

และ อย่างเช่น กล้องฟิล์ม  ก็จะไม่ได้ คนเขมือน แตกต่างแล้วจ้า เป็น กล้องฟิล์ม เปลี่ยนมา เป็น กล้องดิจิทัล

What is Data Mining?

 Data Mining = ข้อมูลที่มีมากมาย สามารถเก็บไว้ แต่มีการ
ผ่านกระบวนการแล้ว

ขั้นตอน

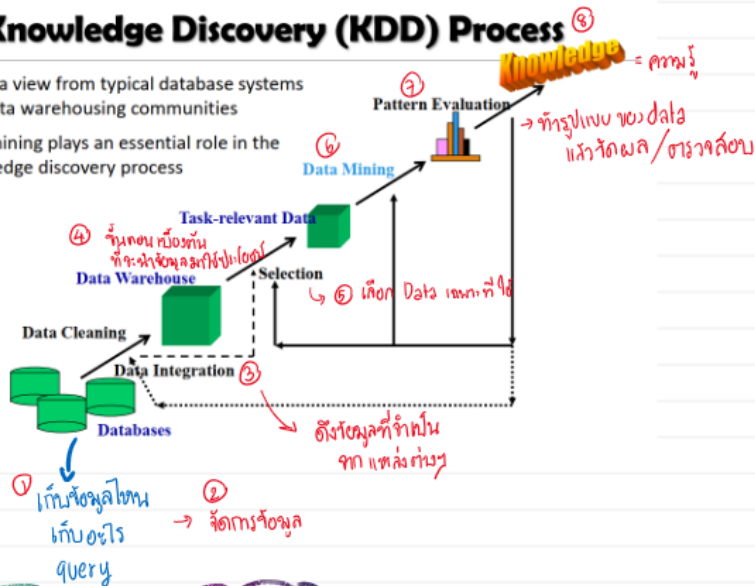
- \rightarrow knowledge discovery in databases "การสกัดองค์ความรู้จาก Database"
- \rightarrow knowledge extraction "การสกัดองค์ความรู้"
- \rightarrow data / pattern analysis "การวิเคราะห์รูปแบบที่มีในข้อมูล"
- \rightarrow data archeology "ค้นคว้าเกี่ยวกับยุคโบราณคดี"
- \rightarrow Business Intelligence "การสกัดองค์ความรู้จาก data เพื่อใช้กับธุรกิจ"



Knowledge Discovery (KDD) Process

Knowledge Discovery (KDD) Process

- This is a view from typical database systems and data warehousing communities
- Data mining plays an essential role in the knowledge discovery process

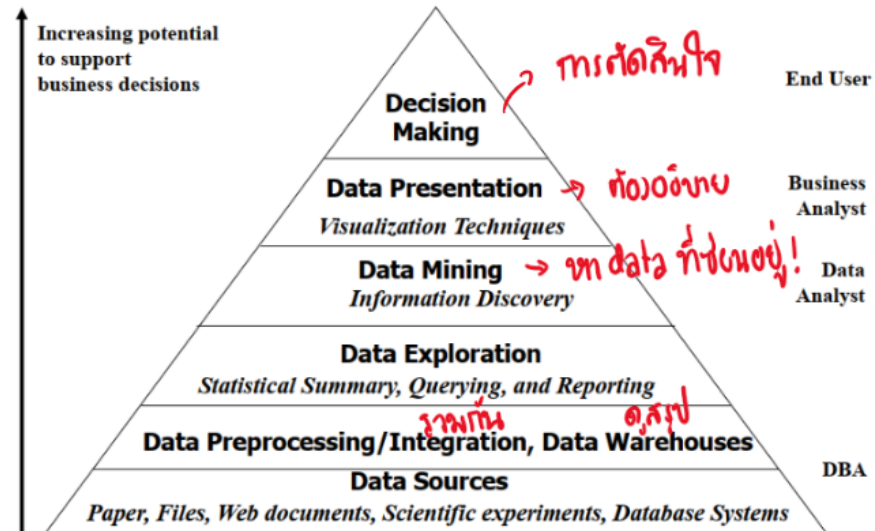


กระบวนการจัดการ Data

Ex มรณ Web mining ?

- Data cleaning
- Data Integration from multiple sources
- Warehousing the data → เป็นข้อมูลที่มีการเก็บไว้ แต่ยังไม่ได้มี การสกัด!
- Data cube construction
- Data selection for data mining
- Data Mining
- Presentation of the mining results —* ต้องนำออกความรู้ มาอธิบายได้
- Patterns & knowledge to be used or stored into knowledge-base

Data Mining in Business Intelligence



นี่คืองานหลักๆ 3 เรื่อง

- "Data Mining"
- ① Pattern discovery = หารูปแบบที่ซ่อนในข้อมูล
 - ② Classification = จำแนก ข้อมูล
 - ③ Clustering = จัดกลุ่มข้อมูล

How the data suppose to look like

Columns (ตาราง) = Attributes, Fields, Features ⇒ ค่าอธิบายคุณสมบัติของข้อมูล

	id	name	domain_id	closed	city_name	zipcode	geohash	new_open	weighted_average_rating	number_of_chains	...	good_for_groups
0	2	Samut Songkhram	2	0	Samut Songkhram	75000	w4rh7g3	0	5.000000	NaN	...	NaN
1	4	Corner House	1	0	Bangkok Metropolitan Region	12150	w4rx73h	0	2.000000	NaN	...	NaN
2	5	ร้านกาแฟ	4	0	Phra Nakhon Si Ayutthaya	13000	w4xd8jk	0	4.000000	NaN	...	NaN
3	6	ร้านกาแฟ	1	0	Bangkok Metropolitan Region	10700.0	w4rqw9q	0	0.000000	NaN	...	NaN
4	7	Buono Caffè	1	0	Bangkok Metropolitan	10220	w4rx4gd	0	3.738462	NaN	...	NaN

row (แถว)

Records, Data point

⇒ ข้อมูลแต่ละตัว

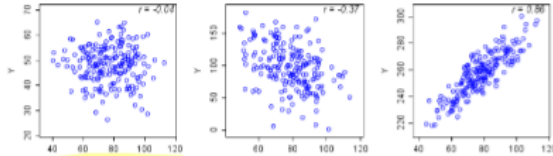
→ 1 จุด

เทคนิคของ Data mining

1

Data Mining Functions: (2) Pattern Discovery

- Frequent patterns (or frequent itemsets)
 - What items are frequently purchased together in your Walmart?
- Association and Correlation Analysis



- A typical association rule —* ทำให้คนรู้จัก Data Mining!
 - Diaper → Beer [0.5%, 75%] (support, confidence)
 - Are strongly associated items also strongly correlated?
 - How to mine such patterns and rules efficiently in large datasets?
 - How to use such patterns for classification, clustering, and other applications?

Association rule

ตัวอย่างที่ทำให้คนรู้จัก => Diaper & Beer

(ตัวอย่าง)

- ⇒ ใช้เทคนิคนี้ในการวิเคราะห์ ใบเสร็จของคนที่มาซื้อของ
- 1 คนที่ซื้อผ้าอ้อม มักจะ ซื้อเบียร์ด้วย
 - 2 และ ก็พบว่า มีคนซื้อผ้าอ้อมกันมากกว่า คู่อื่น ๆ
 - 3 และ มักพบว่า เป็น หากคนพ่อ

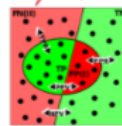
2

Data Mining Functions: (3) Classification

การจำแนกข้อมูล

เลือก 1 Column / 1 attribute เพื่อ เปรียบ Column อื่นๆ

- Classification and label prediction
 - Construct models (functions) based on some training examples
 - Describe and distinguish classes or concepts for future prediction
 - Ex. 1. Classify countries based on (climate) → ฤดูร้อน/ฤดูหนาว
 - Ex. 2. Classify cars based on (gas mileage)
 - Predict some unknown class labels
- Typical methods
 - Decision trees, naïve Bayesian classification, support vector machines, neural networks, rule-based classification, pattern-based classification, logistic regression, ...
- Typical applications:
 - Credit card fraud detection, direct marketing, classifying stars, diseases, web-pages, ...



มีคำตอบ

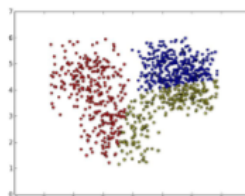
มีค่า y

⇒ เป็นหมวดหมู่

3

Data Mining Functions: (4) Cluster Analysis

- Unsupervised learning (i.e., Class label is unknown)
- Group data to form new categories (i.e., clusters), e.g., cluster houses to find distribution patterns
- Principle: Maximizing intra-class similarity & minimizing interclass similarity
- Many methods and applications



ไม่ใส่คำตอบให้ ทำนาย
แค่ จัดกลุ่ม

ข้อมูลที่มีลักษณะ คล้ายกัน