

สรุปเนื้อหาที่เรียนในคาบที่ 1

เกณฑ์การให้คะแนน

Project Presentation 30%

Programming assignments 40%

Final exam 30%

ดาวน์โหลดเอกสารประกอบการเรียน Data Warehouse and Data mining

จับกลุ่ม กลุ่มละ 5 คน แล้วลงชื่อกลุ่มใน Github และ Excel

สรุปเนื้อหาบทที่ 1

Data mining คือ กระบวนการจัดการกับข้อมูลจำนวนมากเพื่อค้นหารูปแบบและความสัมพันธ์ที่ซ่อนอยู่ในชุดข้อมูลนั้น

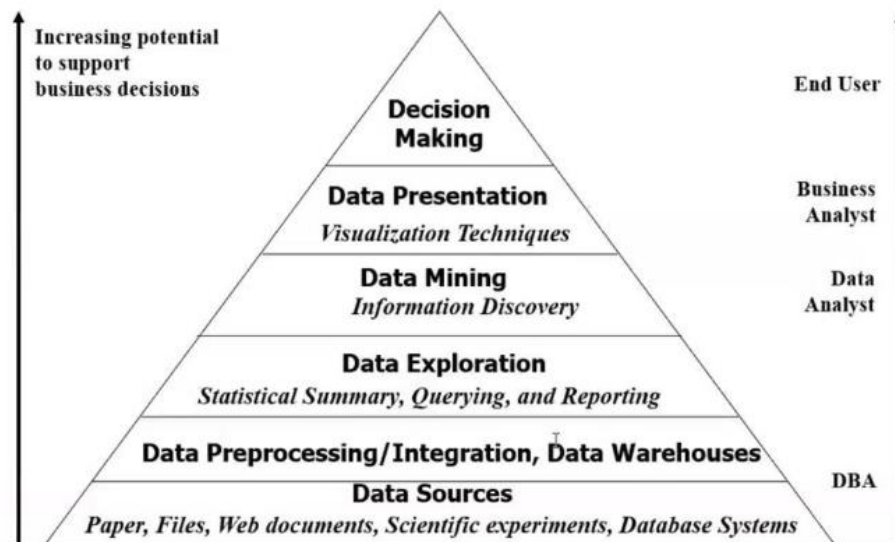
ขั้นตอนการทำ Data mining

ประกอบด้วยขั้นตอนการทำงานย่อยที่จะเปลี่ยนข้อมูลดิบให้กลายเป็นความรู้ ประกอบด้วยขั้นตอนดังนี้

- Data Cleaning เป็นขั้นตอนสำหรับการคัดข้อมูลที่ไม่เกี่ยวข้องออกไป
- Data Integration เป็นขั้นตอนการรวมข้อมูลที่มีหลายแหล่งให้เป็นข้อมูลชุดเดียวกัน
- Data Selection เป็นขั้นตอนการดึงข้อมูลสำหรับการวิเคราะห์จากแหล่งที่บันทึกไว้
- Data Transformation เป็นขั้นตอนการแปลงข้อมูลให้เหมาะสมสำหรับการใช้งาน
- Data Mining เป็นขั้นตอนการค้นหารูปแบบที่เป็นประโยชน์จากข้อมูลที่มีอยู่
- Pattern Evaluation เป็นขั้นตอนการประเมินรูปแบบที่ได้จากการทำเหมืองข้อมูล
- Knowledge Representation เป็นขั้นตอนการนำเสนอความรู้ที่ค้นพบ โดยใช้เทคนิคในการ

นำเสนอเพื่อให้เข้าใจ

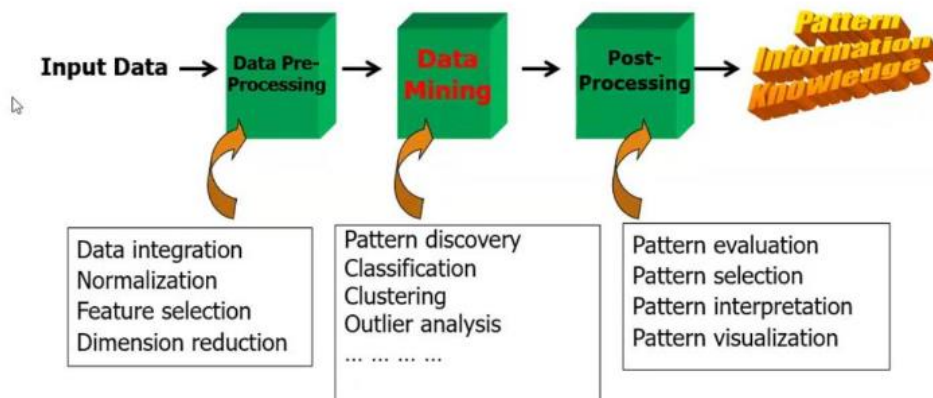
Data Mining in Business Intelligence



12

knowledge discovery (kdd) process

KDD Process: A View from ML and Statistics



□ This is a view from typical machine learning and statistics communities

13

กระบวนการของ Data Mining (A KDD Process เป็นกระบวนการในการค้นหาลักษณะแฝงของข้อมูล (Pattern) ที่ซ่อนอยู่ในฐานข้อมูล

How the data suppose to look like

Columns: Attributes, Fields, Features: ค่าที่ใช้อธิบายคุณสมบัติของข้อมูล

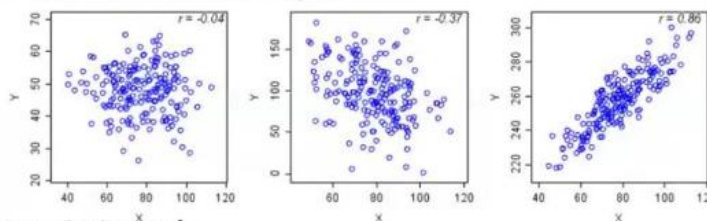
	id	name	domain_id	closed	city_name	zipcode	geohash	new_open	weighted_average_rating	number_of_chains	...	good_for_groups
0	2	นครินทร์ หินตกกรม	2	0	Samut Songkhram	75000	w4rh7g3	0	5.000000	NaN	...	NaN
1	4	Corner House	1	0	Bangkok Metropolitan Region	12150	w4rx73h	0	2.000000	NaN	...	NaN
2	5	วัดโลกยสุธา ราม	4	0	Phra Nakhon Si Ayutthaya	13000	w4x98jk	0	4.000000	NaN	...	NaN
3	6	นันทิดาราโอ กะ	1	0	Bangkok Metropolitan Region	10700.0	w4rqw9q	0	0.000000	NaN	...	NaN
4	7	Buono Caffe	1	0	Bangkok Metropolitan	10220	w4nx4gd	0	3.738462	NaN	...	NaN

Rows: Records, Data point: ข้อมูลแต่ละตัว

เนื้อหาหลักๆที่ได้เรียนมีดังนี้

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
 - ❑ 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 เป็นการค้นหากฎความสัมพันธ์ของข้อมูล โดยค้นหาความสัมพันธ์ของข้อมูลทั้งสองชุดหรือมากกว่าสองชุดขึ้นไปได้ด้วยกัน

Data Mining Functions: (3) Classification

Classification and label prediction

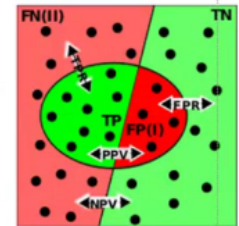
- Construct models (functions) based on some training examples
- Describe and distinguish classes or concepts for future prediction
 - Ex. 1. Classify countries based on (climate)
 - Ex. 2. Classify cars based on (gas mileage)
- Predict some unknown class labels

Typical methods

- Decision trees, naïve Bayesian classification, support vector machines, neural networks, rule-based classification, pattern-based classification, logistic regression, ...

Typical applications:

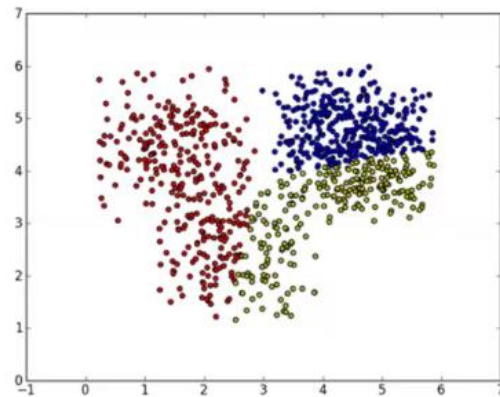
- Credit card fraud detection, direct marketing, classifying stars, diseases, web-pages, ...



Classification เป็นการจัดแบ่งประเภทของข้อมูล โดยหาชุดต้นแบบหรือชุดของการทำงานที่อธิบายและแบ่งประเภทข้อมูล วัตถุประสงค์เพื่อให้สามารถใช้เป็นต้นแบบทำนายประเภทของวัตถุหรือข้อมูลที่ ไม่มีการระบุประเภทหรือชนิดของข้อมูล

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



Clustering คือการจัดกลุ่มข้อมูลซึ่งมีลักษณะคล้ายกับการแบ่งประเภทแต่จะไม่เหมือนกัน โดยการแบ่งประเภทจะวิเคราะห์ข้อมูลตามต้นแบบ แต่สำหรับการแบ่งกลุ่มเป็นการวิเคราะห์โดยไม่พิจารณาจัดกลุ่มตามประเภทที่มีหรือที่รู้จัก แต่จะใช้ขั้นตอนวิธีการจัดกลุ่มเพื่อค้นหากลุ่มที่สามารถยอมรับได้เพื่อจัด เข้ากลุ่ม กล่าวคือ กลุ่มของวัตถุมีการสร้างขึ้นโดยเปรียบเทียบวัตถุที่มีความเหมือนกันจัดเข้า กลุ่มเดียวกัน