



รายงานการวิจัย

SUT Miner: ระบบเหมืองข้อมูลที่มีประสิทธิภาพ (SUT Miner: An efficient data mining system)

ผู้วิจัย

ผู้อำนวยการชุดโครงการวิจัย

รองศาสตราจารย์ ดร.กิตติศักดิ์ เกิดประสพ

สาขาวิชาวิศวกรรมคอมพิวเตอร์

สำนักวิชาวิศวกรรมศาสตร์

หัวหน้าโครงการวิจัยย่อย

รองศาสตราจารย์ ดร. กิตติศักดิ์ เกิดประสพ สาขาวิชาวิศวกรรมคอมพิวเตอร์ สำนักวิชาวิศวกรรมศาสตร์

รองศาสตราจารย์ ดร.นิตยา เกิดประสพ สาขาวิชาวิศวกรรมคอมพิวเตอร์ สำนักวิชาวิศวกรรมศาสตร์

ได้รับทุนอุดหนุนการวิจัยจากมหาวิทยาลัยเทคโนโลยีสุรนารี ปีงบประมาณ พ.ศ. 2548 และ 2549

ผลงานวิจัยเป็นความรับผิดชอบของหัวหน้าโครงการวิจัยแต่เพียงผู้เดียว

มิถุนายน 2553

กิตติกรรมประกาศ

ผู้วิจัยขอขอบคุณมหาวิทยาลัยเทคโนโลยีสุรนารีและสำนักงานคณะกรรมการวิจัยแห่งชาติ ที่ได้จัดสรรงบประมาณในการทำวิจัยให้ในปีงบประมาณ 2548-2549 โครงการนี้ยังได้รับงบประมาณบางส่วน รวมถึงความร่วมมือในการดำเนินงานจากหน่วยปฏิบัติการวิจัยด้านวิศวกรรมข้อมูลและการค้นหาคำรู้ (Data Engineering and Knowledge Discovery -- DEKD -- Research Unit) ขอขอบคุณผู้ทรงคุณวุฒิที่ได้เสียสละเวลาทำหน้าที่ตรวจข้อเสนอโครงการ ตรวจร่างรายงานการวิจัยฉบับสมบูรณ์ และให้ข้อเสนอแนะที่เป็นประโยชน์อย่างมากต่อการปรับปรุงรายงานการวิจัยฉบับสมบูรณ์นี้ ให้ความถูกต้องครบถ้วนมากยิ่งขึ้น

บทคัดย่อภาษาไทย

การทำเหมืองข้อมูลเป็นเทคโนโลยีใหม่ของการวิเคราะห์ข้อมูลอัตโนมัติ เพื่อค้นหาความรู้ที่จะเป็นประโยชน์แก่การวางแผนดำเนินการของหน่วยงานเจ้าของข้อมูล ความรู้ที่ค้นพบนี้เป็นได้หลายรูปแบบ เช่น แพทเทิร์นที่เกิดขึ้นภายในกลุ่มข้อมูลที่สามารถใช้ทำนายลักษณะที่จะเกิดขึ้นในอนาคตหรือทำนายแนวโน้มการเปลี่ยนแปลงของข้อมูล ลักษณะที่สัมพันธ์เชื่อมโยงกันของข้อมูล ลักษณะที่เบี่ยงเบนไปจากข้อมูลกลุ่มใหญ่ และรูปแบบประเภทอื่นๆ อีกหลากหลาย บัจฉยศาสตร์ของการทำเหมืองข้อมูลให้ความรู้ที่มีประโยชน์ คือ อัลกอริทึมสังเคราะห์ความรู้ และ ตัวข้อมูล ถ้าอัลกอริทึมไม่มีประสิทธิภาพก็จะไม่สามารถค้นหาความรู้ที่แฝงอยู่ในข้อมูลได้ หรือถ้าข้อมูลไม่มีคุณภาพเพียงพอก็จะไม่ช่วยให้สังเคราะห์ความรู้ได้ออกมาได้ งานวิจัยนี้จึงเสนอขึ้นเพื่อพัฒนาซอฟต์แวร์ชื่อเอสยูทีไมเนอร์ที่มีส่วนประกอบทั้งส่วนฟิลเตอร์ ส่วนค้นหาแพทเทิร์น และส่วนประเมินแพทเทิร์น เอสยูทีไมเนอร์คือระบบเหมืองข้อมูลที่พัฒนาขึ้น เพื่อให้เป็นเครื่องมือวิเคราะห์ข้อมูลอัจฉริยะสำหรับค้นหารูปแบบหรือสารสนเทศที่เป็นประโยชน์ จากข้อมูลที่เก็บอยู่ในฐานข้อมูล ผู้วิจัยได้นำเสนอวิธีการออกแบบและการพัฒนาเอสยูทีไมเนอร์ ที่จัดเป็นระบบเหมืองข้อมูลสมบูรณ์แบบ เนื่องจากได้ร่วมส่วนการประมวลผลก่อนและหลังการทำเหมืองข้อมูลเข้าไปในระบบเดียวกัน ในรายงานนี้ผู้วิจัยได้นำเสนอกรอบแนวคิดของระบบ และวิธีการพัฒนาระบบในส่วนของโปรแกรมการทำเหมือง โดยใช้วิธีการโปรแกรมเชิงตรรกะด้วยภาษาโปรแกรมมิ่ง วิธีการโปรแกรมเชิงประกาศของโปรแกรมมิ่งช่วยในการเขียนคำสั่งทำได้สั้นและชัดเจน นอกจากนี้การสนับสนุนการทำแพทเทิร์นแมตชิ่งของภาษาโปรแกรมมิ่งยังเป็นข้อได้เปรียบอย่างมากสำหรับงานการค้นหาแบบ ผลสำเร็จของการพัฒนาซอฟต์แวร์นี้จะเป็นจุดเริ่มต้นของการสร้างระบบทำเหมืองข้อมูลสมรรถนะสูงที่ทำงานกับข้อมูลขนาดใหญ่มากและมีลักษณะการเกิดขึ้นอย่างต่อเนื่องได้ต่อไปในอนาคต

บทคัดย่อภาษาอังกฤษ

Data mining is a new technology in automatic data analysis. It is the search for valuable information, or knowledge, in large volumes of data. The discovered knowledge can greatly benefit the organization that owns the data in many ways, for example, to aid decision-making, to reduce the risk of business planning or to project the revenue in future investment. The common types of knowledge discovery are prediction, deviation detection, database segmentation, clustering, association rules and link analysis. Despite the many forms and functions of discovered knowledge, the most important factors in mining valuable knowledge are the learning algorithm and data. The practical data mining system needs an efficient algorithm and high quality data. This research project presents the design and implementation of a complete data mining system, called SUT Miner. The SUT Miner is a data mining system developed as an intelligent data analysis tool to discover patterns and extract useful information from facts stored in databases. In addition to the mining engine, our system incorporates the pre-mining and post-mining parts. In this report, we describe a framework of SUT Miner and present the implementation scheme on the mining engine part using a logic programming paradigm, a Prolog language in particular. A high-level declarative style of Prolog facilitates a clear and concise coding. The language also supports pattern matching, which is a big advantage for a task of pattern discovery. The success of this project will be the preliminary step toward the future development of a high-performance data mining system that can handle very large volume of continuously generated data.

สารบัญ

	หน้า
กิตติกรรมประกาศ	ก
บทคัดย่อภาษาไทย	ข
บทคัดย่อภาษาอังกฤษ	ค
สารบัญ	ง
สารบัญตาราง	ช
สารบัญภาพ	ซ
บทที่ 1 บทนำ	
1.1 ความสำคัญและที่มาของปัญหาที่ทำการวิจัย	1
1.2 วัตถุประสงค์หลักของแผนงานวิจัย	8
1.3 กรอบของแผนงานวิจัย	8
บทที่ 2 วิธีการเตรียมข้อมูลอัตโนมัติก่อนการทำเหมืองข้อมูล (โครงการวิจัยที่ 1)	
2.1 วิธีดำเนินการวิจัยโครงการวิจัยที่ 1	10
2.1.1 กรอบแนวคิดของโครงการวิจัยที่ 1	10
2.1.2 การออกแบบและพัฒนาวิธีการแปลงและปรับปรุงข้อมูล	12
2.1.3 การออกแบบและพัฒนาวิธีการคัดเลือกข้อมูล	16
2.1.4 การคัดเลือกข้อมูลตามความหนาแน่น	17
2.2 การทดสอบโปรแกรมเตรียมข้อมูลอัตโนมัติก่อนการทำเหมืองข้อมูล	19
2.2.1 การทดสอบ Data transformation	19
2.2.2 การทดสอบ Data cleaning	21
2.2.3 การทดสอบ Feature selection	23
2.2.4 การทดสอบ Sampling with replacement	24
2.2.5 การทดสอบ Sampling without replacement	25
2.2.6 การทดสอบ Density-biased sampling	26
2.3 สรุปผลโครงการวิจัยที่ 1	28
บทที่ 3 การพัฒนาการทำเหมืองข้อมูลแบบจัดกลุ่ม (โครงการวิจัยที่ 2)	
3.1 วิธีดำเนินการวิจัยโครงการวิจัยที่ 2	33
3.1.1 กรอบแนวคิดของโครงการวิจัยที่ 2	33
3.1.2 การออกแบบอัลกอริทึมเพื่อจัดกลุ่มข้อมูลตามความหนาแน่น	38
3.1.3 การจัดกลุ่มข้อมูลแบบเพิ่มพูน	42

3.2 การทดสอบโปรแกรมการทำเหมืองข้อมูลแบบจัดกลุ่ม	46
3.2.1 วิธีการทดสอบประสิทธิภาพของการจัดกลุ่ม	46
3.2.2 ผลการทดสอบและอภิปรายผล	47
3.3 สรุปผลโครงการวิจัยที่ 2	53
บทที่ 4 การพัฒนาการทำเหมืองข้อมูลแบบจำแนก (โครงการวิจัยที่ 3)	
4.1 วิธีดำเนินการวิจัยโครงการวิจัยที่ 3	56
4.1.1 กรอบแนวคิดของโครงการวิจัยที่ 3	56
4.1.2 การออกแบบอัลกอริทึมเพื่อการทำเหมืองข้อมูลแบบจำแนก	58
4.1.3 การพัฒนาโปรแกรมเพื่อการทำเหมืองข้อมูลแบบจำแนก	61
4.2 การทดสอบโปรแกรมการทำเหมืองข้อมูลแบบจำแนก	66
4.2.1 วิธีการทดสอบความถูกต้องและประสิทธิภาพของโปรแกรม	66
4.2.2 ผลการทดสอบ	68
4.3 สรุปผลโครงการวิจัยที่ 3	75
บทที่ 5 การประมวลผลหลังกระบวนการทำเหมืองข้อมูล (โครงการวิจัยที่ 4)	
5.1 วิธีดำเนินการวิจัยโครงการวิจัยที่ 4	78
5.1.1 กรอบแนวคิดของโครงการวิจัยที่ 4	78
5.1.2 การออกแบบอัลกอริทึมเพื่อการประเมินและคัดเลือกกฎ	79
5.1.3 การพัฒนาโปรแกรมเพื่อการประมวลผลหลังการทำเหมืองข้อมูล ...	82
5.2 การทดสอบโปรแกรมเพื่อการประมวลผลหลังการทำเหมืองข้อมูล	92
5.2.1 วิธีการทดสอบโปรแกรมเพื่อการประมวลผลหลังการทำเหมือง ข้อมูล	92
5.2.2 ผลการทดสอบโปรแกรม	96
5.3 สรุปผลโครงการวิจัยที่ 4	99
บทที่ 6 บทสรุป	
6.1 ลักษณะเด่นของระบบ SUT Miner	104
6.2 ข้อจำกัดของระบบ SUT Miner	105
6.3 ข้อเสนอแนะและแนวทางการพัฒนาในอนาคต	106
บรรณานุกรม	108

ภาคผนวก	
ภาคผนวก ก คู่มือการใช้งานระบบ SUT Miner	113
ภาคผนวก ข รหัสต้นฉบับของชุดโปรแกรม SUT Miner	121
ภาคผนวก ค ผลงานวิจัยของชุดโครงการ SUT Miner ที่ได้รับการตีพิมพ์เผยแพร่	146
ประวัติผู้วิจัย	205

สารบัญตาราง

	หน้า
ตารางที่ 3.1 เปรียบเทียบเวลาและหน่วยความจำที่ใช้ในแต่ละอัลกอริทึม	36
ตารางที่ 3.2 เปรียบเทียบค่า means ของการจัดกลุ่มแบบ batch และแบบ incremental	45
ตารางที่ 3.3 ผลการจัดข้อมูลเข้ากลุ่มของวิธีการจัดกลุ่มแบบ batch และแบบ incremental	46
ตารางที่ 3.4 เปรียบเทียบผลการจัดกลุ่มข้อมูล post-operative ด้วยวิธีจัดกลุ่มแบบ batch และแบบ incremental	49
ตารางที่ 3.5 เปรียบเทียบผลการจัดกลุ่มข้อมูล breast cancer ด้วยวิธีจัดกลุ่มแบบ batch และแบบ incremental	50
ตารางที่ 3.6 ค่า SSE ของการจัดกลุ่มข้อมูลแบบ batch และแบบ incremental ของข้อมูล post-operative	51
ตารางที่ 3.7 ค่า SSE ของการจัดกลุ่มข้อมูลแบบ batch และแบบ incremental ของข้อมูล breast cancer	52
ตารางที่ 4.1 รายละเอียดของข้อมูลที่ใช้ทดสอบความถูกต้องและประสิทธิภาพของ โมเดล	68
ตารางที่ 4.2 ผลการทดสอบความแม่นยำของโมเดล Probabilistic decision rules เปรียบเทียบกับ โมเดล ID3	68
ตารางที่ 4.3 ความแม่นยำของโมเดล Probabilistic decision rules เมื่อมีการลดขนาด โมเดล	70
ตารางที่ 5.1 ผลการทดสอบความถูกต้องของระบบผู้เชี่ยวชาญกับข้อมูล post-operative patients	96
ตารางที่ 5.2 ผลการทดสอบความถูกต้องของระบบผู้เชี่ยวชาญกับข้อมูล breast-cancer recurrences	97
ตารางที่ 5.3 ผลการวิเคราะห์ความผิดพลาดในลักษณะของ false negative	99

สารบัญญภาพ

	หน้า
รูปที่ 1.1 กระบวนการทำเหมืองข้อมูลเพื่อคัดแยกความรู้ออกจากข้อมูลดิบ	3
รูปที่ 1.2 กระบวนการทำเหมืองข้อมูลโดยสรุป	4
รูปที่ 1.3 แผนภาพรวมของระบบเหมืองข้อมูล SUT Miner	9
รูปที่ 2.1 โครงสร้างซอฟต์แวร์เตรียมข้อมูลก่อนการทำเหมืองข้อมูล	10
รูปที่ 2.2 ตัวอย่างไฟล์ข้อมูลในรูปแบบ UCI repository	13
รูปที่ 2.3 ตัวอย่างไฟล์ข้อมูลในรูปแบบ Horn clauses	13
รูปที่ 2.4 ตัวอย่างไฟล์ข้อมูลที่ปรากฏ missing values	15
รูปที่ 2.5 แนวคิดของการสุ่มข้อมูลตามความหนาแน่น	18
รูปที่ 2.6 ข้อมูลที่ใช้ทดสอบการทำงานของโปรแกรมเตรียมข้อมูล	20
รูปที่ 2.7 จอภาพส่วน Data transformation	20
รูปที่ 2.8 จอภาพของ SWI Prolog ขณะแปลงข้อมูลและไฟล์ golf_out ที่ได้	21
รูปที่ 2.9 ข้อมูลในไฟล์ golf_out ที่ค่าที่สูญหายถูกแทนที่ด้วยข้อความ 'missing'	22
รูปที่ 2.10 ข้อมูลในไฟล์ golf_out ที่ค่าที่สูญหายถูกแทนที่ด้วยค่าส่วนใหญ่	22
รูปที่ 2.11 จอภาพส่วน Feature selection และ Data sampling	23
รูปที่ 2.12 ผลลัพธ์ของ Feature selection โดยระบุแอททริบิวต์ [outlook,humidity,windy]	23
รูปที่ 2.13 ผลลัพธ์ของการสุ่มข้อมูลแบบมีการใส่ค่ากลับคืน	24
รูปที่ 2.14 การสุ่มข้อมูลที่ขนาด 50% และ ไม่มีการใส่ค่ากลับคืน	25
รูปที่ 2.15 การสุ่มตามความหนาแน่นด้วยเกณฑ์ขั้นต่ำ [3, 0.23]	26
รูปที่ 2.16 ผลลัพธ์ของการสุ่มตามความหนาแน่นด้วยเกณฑ์ขั้นต่ำ [1, 0.14]	27
รูปที่ 2.17 พัฒนาการของข้อมูลที่ถูกเปลี่ยนเป็นสารสนเทศและความรู้	28
รูปที่ 2.18 ขั้นตอนต่างๆ ในกระบวนการทำเหมืองข้อมูล	29
รูปที่ 3.1 แสดงการจัดหมวดหมู่เทคนิคที่ใช้ในการทำ clustering	33
รูปที่ 3.2 การจัดกลุ่มข้อมูลด้วยเทคนิค partitioning	34
รูปที่ 3.3 การจัดกลุ่มข้อมูลด้วยเทคนิค hierarchical	35
รูปที่ 3.4 กรอบแนวคิดของงานออกแบบและพัฒนาการจัดกลุ่มข้อมูลตามความหนาแน่น	37
รูปที่ 3.5 ตัวอย่างไฟล์ข้อมูลที่จะนำเข้ายัง โปรแกรมจัดกลุ่มข้อมูลตามความหนาแน่น	39
รูปที่ 3.6 จอภาพเริ่มต้นของ โปรแกรมจัดกลุ่มข้อมูลตามความหนาแน่น	40
รูปที่ 3.7 ผลลัพธ์ของการจัดกลุ่มข้อมูลเมื่อระบุค่าความหนาแน่นเป็นศูนย์	40
รูปที่ 3.8 ผลลัพธ์ของการจัดกลุ่มข้อมูลเมื่อระบุค่าความหนาแน่นเป็น [2, 0.143]	41
รูปที่ 3.9 ผลลัพธ์ของการจัดข้อมูลเป็นสามกลุ่มและระบุค่าความหนาแน่นเป็น [3, 0.143]	42

รูปที่ 3.10	จอภาพของโปรแกรมการรวมคลัสเตอร์ในงานจัดกลุ่มข้อมูลแบบเพิ่มพูน	43
รูปที่ 3.11	ผลลัพธ์ของการจัดกลุ่มข้อมูลแบบเพิ่มพูนกับข้อมูล weather	44
รูปที่ 3.12	การทำ unsupervised learning ด้วย clustering algorithm	54
รูปที่ 3.13	การทำ supervised learning ด้วย classification algorithm	54
รูปที่ 4.1	กรอบแนวคิดของการออกแบบและพัฒนาซอฟต์แวร์เพื่อการทำเหมืองข้อมูล แบบจําแนก	56
รูปที่ 4.2	ตัวอย่างไฟล์ข้อมูลที่จะนำเข้ายังโปรแกรมเพื่อสร้างกฎการตัดสินใจ	61
รูปที่ 4.3	หน้าจอหลักของการเรียกใช้โปรแกรมสร้างต้นไม้ตัดสินใจ	62
รูปที่ 4.4	รายละเอียดที่แสดงจากการใช้คำสั่ง listing(node) และ listing(edge)	63
รูปที่ 4.5	โมเดลที่ได้ในลักษณะของกฎการตัดสินใจเมื่อเรียกใช้โปรแกรมด้วยค่า ความน่าจะเป็น 0.0	65
รูปที่ 4.6	โมเดลที่ได้เมื่อเรียกใช้โปรแกรมด้วยค่าความน่าจะเป็น 0.155	65
รูปที่ 4.7	กราฟเปรียบเทียบความแม่นยำ (accuracy) ของ Probabilistic decision rules และ ID3	69
รูปที่ 4.8	ความแม่นยำของโมเดล Probabilistic decision rules ที่ขนาดต่างๆ กันของชุด ข้อมูล Monk	71
รูปที่ 4.9	ความแม่นยำของโมเดล Probabilistic decision rules ที่ขนาดต่างๆ กันของชุด ข้อมูล Post-operative	71
รูปที่ 4.10	ความแม่นยำของโมเดล Probabilistic decision rules ที่ขนาดต่างๆ กันของชุด ข้อมูล Breast cancer	72
รูปที่ 4.11	ความแม่นยำของโมเดล Probabilistic decision rules ที่ขนาดต่างๆ กันของชุด ข้อมูล Vote	72
รูปที่ 4.12	ความแม่นยำของโมเดล Probabilistic decision rules ที่ขนาดต่างๆ กันของชุด ข้อมูล Hepatitis	73
รูปที่ 4.13	ความแม่นยำของโมเดล Probabilistic decision rules ที่ขนาดต่างๆ กันของชุด ข้อมูล Mushroom	73
รูปที่ 5.1	สถาปัตยกรรมของระบบจัดการความรู้ที่ได้จากการทำเหมืองข้อมูล	78
รูปที่ 5.2	ข้อมูลที่ใช้เพื่อสร้างกฎการตัดสินใจ	83
รูปที่ 5.3	จอภาพให้ผู้ใช้ระบุชื่อข้อมูลและโมเดลข้อมูลในลักษณะ node และ edge	85
รูปที่ 5.4	ข้อมูลในไฟล์ 1.knb	87

รูปที่ 5.5 โครงสร้างของ expert system shell	87
รูปที่ 5.6 ระบบผู้เชี่ยวชาญที่ให้คำแนะนำเกี่ยวกับคอนแทกเลนส์	91
รูปที่ 5.7 การโต้ตอบของระบบผู้เชี่ยวชาญกรณีไม่มีข้อมูลปรากฏในฐานความรู้	91
รูปที่ 5.8 ตัวอย่างการทดสอบความถูกต้องของระบบผู้เชี่ยวชาญ	93
รูปที่ 5.9 กฎทั้งหมดของระบบผู้เชี่ยวชาญที่สร้างด้วยเกณฑ์ความน่าจะเป็นขั้นต่ำ 0.001	94
รูปที่ 5.10 ฐานความรู้ของระบบผู้เชี่ยวชาญในการแนะนำเกี่ยวกับ breast-cancer recurrences ...	95
รูปที่ 5.11 กราฟเปรียบเทียบ error rate เมื่อทดสอบกับข้อมูล post-operative patients	98
รูปที่ 5.12 กราฟเปรียบเทียบ error rate เมื่อทดสอบกับข้อมูล breast-cancer recurrences	98
รูปที่ 5.13 โมเดลของข้อมูล post-operative patients สร้างจากโปรแกรม ID3	100
รูปที่ 5.14 โมเดลของข้อมูล post-operative patients สร้างจากโปรแกรม Multi-layer perceptron	101
รูปที่ 5.15 โมเดลของข้อมูล post-operative patients ที่มีค่าความน่าจะเป็นสูงกว่า 0.02.....	102
รูปที่ 5.16 การใช้ expert system shell สอบถาม โมเดลข้อมูล	102
รูปที่ 6.1 ระบบเหมืองข้อมูล SUT Miner และ โมดูลหลักในแต่ละส่วน	104
รูปที่ ก1 การใช้คำสั่งเพื่อนำเข้าเพิ่มข้อมูล	114
รูปที่ ก2 การระบุชื่อเพิ่มข้อมูลเข้าและไฟล์ใหม่ที่เกิดขึ้นหลังการแปลงรูปแบบข้อมูล	114
รูปที่ ก3 จอภาพส่วนคัดเลือกแอททริบิวต์และสุ่มข้อมูล	115
รูปที่ ก4 จอภาพเริ่มต้นของ โปรแกรมจัดกลุ่มข้อมูล	116
รูปที่ ก5 ผลลัพธ์ของการจัดกลุ่มข้อมูลเมื่อระบุจำนวนกลุ่มเป็นสองกลุ่ม	116
รูปที่ ก6 การจัดกลุ่มข้อมูลแบบเพิ่มพูนกับข้อมูลในสองกลุ่มย่อย	117
รูปที่ ก7 โมเดลข้อมูลในลักษณะต้นไม้ตัดสินใจเมื่อเรียกใช้โปรแกรมทำเหมืองข้อมูล แบบจำแนก	118
รูปที่ ก8 ฐานความรู้ที่ได้จากการแปลงโมเดลในลักษณะต้นไม้ตัดสินใจ	119
รูปที่ ก9 ตัวอย่างการเรียกดูโมเดลข้อมูลผ่านระบบผู้เชี่ยวชาญ	120

บทที่ 1

บทนำ

1.1 ความสำคัญและที่มาของปัญหาที่ทำการวิจัย

การขุดค้นและวิเคราะห์ข้อมูลในฐานข้อมูลเพื่อค้นหาความรู้ เรียกว่า *การทำเหมืองข้อมูล* (data mining) เป็นเทคโนโลยีใหม่ของการประยุกต์ใช้ข้อมูลที่เกิดขึ้นในฐานข้อมูลให้เกิดประโยชน์สูงสุดแก่หน่วยงานที่เป็นเจ้าของข้อมูล การประยุกต์ใช้ข้อมูลที่กำลังมีได้หลายแนวทาง แต่โดยทั่วไปมักจะเป็นการสรุปภาพรวมของข้อมูลในฐานข้อมูล, การจัดกลุ่มข้อมูลอัตโนมัติ, การวิเคราะห์แนวโน้มการเปลี่ยนแปลงของข้อมูล, หรือ การค้นหาความสัมพันธ์ที่ซ่อนอยู่ภายในกลุ่มของข้อมูล เช่น จากข้อมูลการซื้อสินค้าของลูกค้าซูเปอร์มาเก็ตในรอบหนึ่งเดือนที่ผ่านมา เมื่อวิเคราะห์หาความสัมพันธ์แล้วพบว่า “ถ้าลูกค้าซื้อนมและขนมปังแล้ว ลูกค้าจะซื้อเบียร์สำเร็จรูปด้วย” ผลลัพธ์ที่ได้ก็คือ ความรู้เกี่ยวกับพฤติกรรมของผู้บริโภค ที่จะช่วยให้ผู้จัดการซูเปอร์มาเก็ตสามารถวางแผนการจัดวางสินค้าและการจัดโปรแกรมกระตุ้นยอดขายสินค้าบางรายการได้อย่างมีประสิทธิภาพ

การวิเคราะห์ข้อมูลด้วยโปรแกรมช่วยงาน เช่น SPSS, SAS นักวิเคราะห์ข้อมูลจะต้องเป็นผู้กำหนดว่าจะศึกษาลักษณะใดจากข้อมูลและจะใช้ข้อมูลส่วนใดบ้าง แต่โปรแกรมทำเหมืองข้อมูลจะกระทำขั้นตอนต่างๆ เหล่านี้ให้โดยอัตโนมัติ โปรแกรมทำเหมืองข้อมูลมีความสามารถที่จะค้นหาแนวโน้ม รูปแบบร่วม หรือลักษณะอื่นๆที่น่าสนใจ โดยไม่ต้องพึ่งพาการสั่งงานทุกขั้นตอนจากนักวิเคราะห์ข้อมูล นอกจากนี้ยังสามารถค้นพบลักษณะที่น่าสนใจจากข้อมูลซึ่งนักวิเคราะห์ข้อมูลไม่ได้คาดหมายมาก่อน

ระบบเหมืองข้อมูล (data mining system) มีความแตกต่างจากระบบผู้เชี่ยวชาญ (expert system) ตรงที่ฐานความรู้ของระบบเหมืองข้อมูลได้จากการตั้งเคราะห์จากข้อมูลโดยตรง สามารถปรับปรุงฐานความรู้ของตัวเองได้อัตโนมัติตามข้อมูลใหม่ที่ได้รับเพิ่มขึ้น ซึ่งจะแตกต่างจากระบบผู้เชี่ยวชาญที่ฐานความรู้ถูกป้อนเข้ามาในระบบโดยวิศวกรความรู้ (knowledge engineer) และจะคงตัวอยู่เช่นนั้นตลอดการใช้งานจนกว่าวิศวกรความรู้จะปรับปรุงฐานความรู้

การประยุกต์ใช้เทคโนโลยีการทำเหมืองข้อมูลมีได้หลากหลายลักษณะ สามารถจัดกลุ่มกว้างๆได้เป็นสองกลุ่ม คือ กลุ่มที่ใช้การทำเหมืองข้อมูลเพื่อการทำนาย และกลุ่มที่ใช้เพื่อการอธิบาย

การทำการทำเหมืองข้อมูลเพื่อการทำนาย เป็นการนำความรู้ที่เรียนรู้มาจากข้อมูลที่มีอยู่เพื่อประโยชน์ในการทำนายข้อมูลใหม่ที่จะเกิดขึ้นในอนาคต เช่น จากข้อมูลลูกค้าของแผนกสินค้า

ของธนาคารที่ได้มีการจัดกลุ่มของลูกค้าไว้แล้วว่าใครเป็นลูกค้าชั้นดี ใครเป็นลูกค้าในระดับปานกลาง และใครเป็นลูกค้าที่มักจะผิดนัดชำระหนี้ โปรแกรมการทำเหมืองข้อมูลสามารถเรียนรู้จากข้อมูลเหล่านี้และค้นหาโมเดลที่สามารถใช้อธิบายลักษณะของลูกค้าชั้นดี ลูกค้าระดับปานกลาง และลูกค้าที่ไม่เป็นที่ต้องการ จากโมเดลที่ได้นี้สามารถนำไป ใช้ทำนายลูกค้าใหม่ที่มาขอสินเชื่อได้ ว่าเขาน่าจะเป็นลูกค้าประเภทใด

การทำเหมืองข้อมูลเพื่อการอธิบาย เป็นการค้นหารูปแบบหรือแพทเทิร์นบางอย่างที่น่าสนใจจากกลุ่มข้อมูล รูปแบบนี้มักจะเป็นความสัมพันธ์ หรือลักษณะที่เชื่อมโยงกันของข้อมูล การทำเหมืองข้อมูลแบบนี้ต่างจากแบบแรกตรงที่ ผู้ใช้ไม่ได้กำหนดล่วงหน้าว่าจะให้โปรแกรมทำเหมืองข้อมูลค้นหารูปแบบหรือ โมเดลของอะไร แต่ให้ค้นหาทุกรูปแบบที่น่าสนใจจากข้อมูล

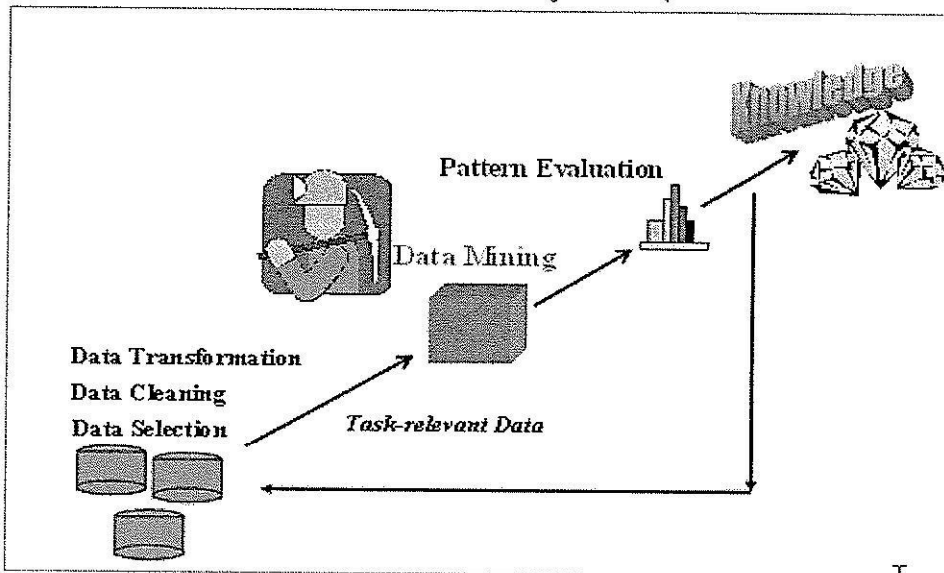
การทำเหมืองข้อมูลทั้งแบบเพื่อการทำนายและเพื่อการอธิบาย จัดเป็นการเรียนรู้ของเครื่องแบบที่ต้องรับการชี้แนะ (supervised learning) การเรียนรู้ของเครื่องในแบบที่ยากขึ้น คือ แบบที่ไม่ต้องชี้แนะ (unsupervised learning) ได้แก่ การจัดกลุ่มข้อมูลอัตโนมัติ (clustering)

งานทำเหมืองข้อมูลที่มีได้หลากหลายรูปแบบนี้ เกิดขึ้นเนื่องจากความรู้ที่ต้องการได้จากข้อมูลมีได้หลายลักษณะ ตัวอย่างต่อไปนี้แสดงผลสำเร็จส่วนหนึ่งของการนำเทคโนโลยีการทำเหมืองข้อมูลไปประยุกต์ใช้

- **ด้านการแพทย์:** ใช้โปรแกรมทำเหมืองข้อมูลค้นหาผลข้างเคียงที่ไม่ได้คาดหมายมาก่อนของการใช้ยา โดยอาศัยข้อมูลจากแฟ้มประวัติผู้ป่วย, ใช้ในการวิเคราะห์หาความสัมพันธ์ของสารพันธุกรรม
- **ด้านการเงิน:** ใช้โปรแกรมทำเหมืองข้อมูลช่วยในการตัดสินใจว่าควรจะอนุมัติเครดิตให้ลูกค้ารายใดบ้าง, ใช้ในการคาดหมายความน่าจะเป็นว่าธุรกิจนั้นๆมีโอกาที่จะล้มละลายหรือไม่, ใช้คาดหมายการขึ้น/ลงของหุ้นในตลาดหลักทรัพย์
- **ด้านการเกษตร:** ใช้โปรแกรมทำเหมืองข้อมูลจำแนกประเภทของโรคพืชที่เกิดกับถั่วเหลืองและมะเขือเทศ
- **ด้านวิศวกรรม:** ใช้โปรแกรมทำเหมืองข้อมูลวิเคราะห์และวินิจฉัยสาเหตุการทำงานผิดพลาดของเครื่องจักรกล
- **ด้านอาชญาวิทยา:** ใช้โปรแกรมทำเหมืองข้อมูลวิเคราะห์หาเจ้าของลายนิ้วมือ
- **ด้านอวกาศ:** ใช้โปรแกรมทำเหมืองข้อมูลวิเคราะห์ข้อมูลที่ส่งมาจากดาวเทียมขององค์การนาซา

กระบวนการทำเหมืองข้อมูล เปรียบได้กับการทำเหมืองแร่ที่เราใช้เครื่องจักรคัดแยกแร่ที่มีค่าออกจากกองหิน กรวด ดินที่ปะปนมากับสายแร่ เพียงแต่ในกระบวนการทำเหมืองข้อมูลสิ่งที่เราได้จากกองข้อมูลมหาศาล คือ ความรู้ (knowledge) ที่ซ่อนอยู่ในกองข้อมูล กระบวนการนี้แสดงเป็นแผนภาพได้ดังรูปที่ 1.1

ความรู้ที่ได้นี้จะช่วยให้เราเข้าใจลักษณะของข้อมูล และเข้าใจปัจจัยที่ทำให้เกิดลักษณะบางอย่างขึ้นในข้อมูลบางกลุ่ม ซึ่งจะช่วยให้เราสามารถทำนายแนวโน้มของข้อมูลใหม่ที่จะเกิดขึ้นในอนาคตได้ รวมถึงเข้าใจความสัมพันธ์ที่เชื่อมโยงข้อมูลแต่ละกลุ่มย่อยเข้าด้วยกัน



รูปที่ 1.1 กระบวนการทำเหมืองข้อมูลเพื่อคัดแยกความรู้ออกจากข้อมูลดิบ

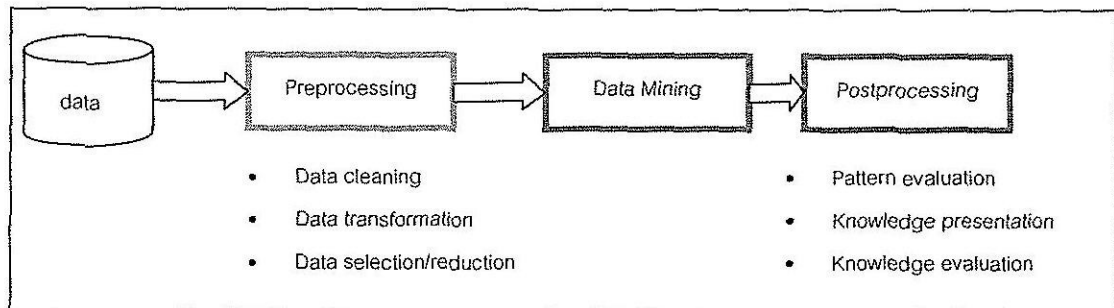
โดยทั่วไปกระบวนการทำเหมืองข้อมูล จะประกอบด้วย 4 ขั้นตอนใหญ่ๆคือ

- 1) เตรียมข้อมูล (data preparation) : ประกอบด้วยขั้นตอนย่อย คือ การคัดเลือกข้อมูล (data selection) เพื่อให้ได้เฉพาะข้อมูลที่จะเป็นประโยชน์ต่อการทำเหมืองข้อมูล การคัดเลือกข้อมูลอย่างเหมาะสม จะช่วยให้ค้นหาความรู้ที่เป็นประโยชน์ได้อย่างรวดเร็ว นอกจากนี้ยังต้องมีการปรับปรุงคุณภาพข้อมูล (data cleaning) ซึ่งจำเป็นในกรณีที่ข้อมูลดิบมีส่วนที่ไม่สมบูรณ์ หรือมีข้อมูลบางส่วนบกพร่อง มีข้อมูลรบกวน (noise)ปะปนอยู่ และถ้าข้อมูลไม่อยู่ในรูปแบบที่ถูกต้องหรือเหมาะสม จะต้องมีการปรับเปลี่ยนรูปแบบข้อมูล (data transformation) ให้อยู่ในรูปแบบที่โปรแกรมทำเหมืองข้อมูลจะเรียกใช้งานได้
- 2) ค้นหาโมเดลจากข้อมูล (mining) : กระบวนการค้นหาโมเดลหรือความสัมพันธ์ จะเริ่มจากข้อมูลเริ่มต้นจำนวนไม่มากนัก จากนั้นนำผลที่ได้จากอัลกอริทึมค้นหาโมเดล (mining algorithm) ไปยืนยันกับข้อมูลทดสอบ ถ้าผลที่ได้ยังไม่น่าพอใจอาจจะต้อง

ปรับค่าพารามิเตอร์บางตัวของ mining algorithm และเริ่มกระบวนการค้นหาใหม่กับข้อมูลจำนวนมากขึ้น จนกว่าผลที่ได้มีความถูกต้องอยู่ในระดับที่ยอมรับได้ จึงจะจบกระบวนการค้นหา

- 3) ตรวจสอบและวิเคราะห์ผล (pattern evaluation) : โมเดลหรือความสัมพันธ์ที่หามาได้ในขั้นตอนที่ 2 จะต้องถูกนำมาทดสอบอัตราความผิดพลาดและวิเคราะห์ความซับซ้อนของรูปแบบโมเดล ถ้าอัตราความผิดพลาดยังสูงเกินไป อาจจะต้องย้อนกลับไปขั้นตอนที่ 2 อีกครั้ง เพื่อปรับปรุงโมเดลให้ถูกต้องยิ่งขึ้น ในทำนองเดียวกัน ถ้าโมเดลที่หามาได้มีรูปแบบที่ซับซ้อนเกินไปจนยากต่อการทำความเข้าใจ อาจจะต้องย้อนกระบวนการกลับไป เพื่อให้หาโมเดลใหม่ที่มีความถูกต้องเท่าเดิมแต่มีรูปแบบที่ซับซ้อนน้อยลง
- 4) กัดเลือกโมเดลและกำหนดวิธีแสดงผล (knowledge processing and representation) : โมเดลหรือความรู้ที่ได้จาก mining algorithm มักจะมีปริมาณมาก และรูปแบบการแสดงความรู้ที่ยากต่อการทำความเข้าใจ จึงต้องมีกระบวนการภายหลังการทำเหมืองข้อมูล (postprocessing) เพื่อประเมินและคัดเลือกความรู้ที่น่าสนใจ

จากทั้งสี่ขั้นตอนในกระบวนการทำเหมืองข้อมูล เริ่มต้นตั้งแต่การนำเข้าสู่ข้อมูลจนกระทั่งได้ความรู้ที่นำไปใช้ประโยชน์ได้ สรุปภาพรวมของกระบวนการได้ดังรูปที่ 1.2



รูปที่ 1.2 กระบวนการทำเหมืองข้อมูล โดยสรุป

จากภาพรวมของกระบวนการทำเหมืองข้อมูล จะเห็นได้ว่าการทำเหมืองข้อมูลเป็นศาสตร์ที่ต้องรวมความก้าวหน้าจากหลายสาขา ได้แก่ database technology, artificial intelligence, machine learning, statistics, neural networks, pattern recognition, knowledge-based systems, knowledge acquisition, information retrieval, high-performance computing และ data visualization การทำเหมืองข้อมูลเป็นศาสตร์แขนงใหม่ที่เริ่มเป็นที่รู้จักในช่วงปลายทศวรรษที่ 1980 Piatetsky-Shapiro และ Frawley (1991) ได้รวบรวมงานวิจัยในยุคแรกของศาสตร์นี้ไว้ในหนังสือชื่อ

Knowledge Discovery in Databases และได้เริ่มมีการรวมตัวของนักวิจัยในสาขานี้จัดการประชุมทางวิชาการเพื่อเสนอความก้าวหน้าของงานวิจัยภายใต้ชื่อการประชุม International Conference on Knowledge Discovery and Data Mining (KDD) ตั้งแต่ปี 1995 ปัจจุบันศาสตร์ทางด้านการทำเหมืองข้อมูลได้รับความสนใจจากนักวิจัยจำนวนมากทั่วโลกในสาขาต่างๆ ที่เกี่ยวข้อง ค้นคว้าเพื่อหาแนวทางที่จะทำให้กระบวนการทำเหมืองข้อมูลในทุกขั้นตอนมีประสิทธิภาพสูงขึ้น และทำงานกับข้อมูลขนาดใหญ่มาได้

การทำเหมืองข้อมูลให้มีประสิทธิภาพ จะต้องเริ่มตั้งแต่การเตรียมข้อมูลให้ได้ข้อมูลที่มีคุณภาพ นักวิจัยจำนวนมาก เช่น Redman (1992), Wang et al. (1995), Wand และ Wang (1996), Kennedy et al. (1998), Weiss และ Indurkha (1998), Pyle (1999), Ballou และ Tayi (1999) ได้อธิบายลักษณะของข้อมูลที่มีคุณภาพและเสนอแนะเทคนิคที่สามารถนำมาใช้ช่วยเพิ่มคุณภาพข้อมูลลักษณะที่มักจะพบบ่อยในกลุ่มของข้อมูลที่ไม่มีคุณภาพ ได้แก่ ข้อมูลบางส่วนขาดหายไป (missing values) มีผู้เสนอแนวทางจัดการกับกรณีเช่นนี้ไว้หลายแนวทางดังปรากฏในงานวิจัยของ Friedman (1977), Breiman et al. (1984) และ Quinlan (1989)

ประโยชน์ของการทำเหมืองข้อมูลจะเห็นได้ชัดเจนเมื่อข้อมูลมีปริมาณมาก แต่ปัญหาที่เกิดขึ้นตามมาคือระบบเหมืองข้อมูลไม่สามารถรองรับข้อมูลขนาดมหาศาลนั้นได้ แนวทางแก้ปัญหานี้คือใช้วิธีลดขนาดข้อมูล เทคนิคการลดขนาดข้อมูลมักจะกระทำในสองแนวทาง คือ ลดลักษณะ (attributes or features) ที่อธิบายข้อมูลแต่ละรายการให้เหลือเฉพาะลักษณะหลักที่สำคัญ งานวิจัยในแนวทางนี้จะอยู่ในกลุ่มที่เรียกว่า feature selection รายละเอียดปรากฏในหนังสือและงานวิจัยจำนวนมาก (Neter et al., 1996; Dash and Liu, 1997; Kohavi and John, 1997; Dash et al., 1997; Liu and Motoda, 1998)

แนวทางที่สองในการลดขนาดของข้อมูล ใช้วิธีการลดจำนวนรายการข้อมูลให้เหลือเฉพาะข้อมูลที่สามารถเป็นตัวแทนของข้อมูลกลุ่มใหญ่ได้ เทคนิคที่นิยมใช้ลดจำนวนข้อมูลคือ การสุ่ม (sampling) การจัดกลุ่ม (clustering) และการใช้ฮิสโตแกรม (histogram) John และ Langley (1996) ได้ศึกษา static และ dynamic sampling ในขณะที่ Josien และคณะ (2001) ได้ทดลองศึกษาพฤติกรรมของโปรแกรมทำเหมืองข้อมูลเมื่อข้อมูลถูกลดขนาดลง Kivinen และ Mannila (1993) ได้วิเคราะห์ในเชิงทฤษฎีถึงจำนวนข้อมูลที่สามารถเป็นตัวแทนของข้อมูลกลุ่มใหญ่ได้ Barbara และคณะ (1997) รวมถึง Devore และ Peck (1997) ได้นำเทคนิคฮิสโตแกรมมาใช้วิเคราะห์การกระจายของข้อมูลเพื่อสุ่มข้อมูลจากแต่ละกลุ่ม

งานวิจัยและพัฒนาทั้งหลายที่เกี่ยวข้องกับการเตรียมข้อมูลมักจะทำเพียงเฉพาะบางส่วนของ การเตรียมข้อมูล เช่น แก้ไขในกรณีข้อมูลบางส่วนสูญหาย (missing values) หรือเสนอเทคนิคเฉพาะบางเทคนิคเพื่อให้ข้อมูลมีขนาดที่เหมาะสมกับโปรแกรมทำเหมืองข้อมูล แต่โครงการวิจัยนี้

ได้พัฒนาทั้งสามขั้นตอนของการเตรียมข้อมูลคือ data transformation, data cleaning และ data reduction ข้อมูลที่ผ่านการเตรียมแล้วจะสามารถถูกส่งต่อไปให้กับระบบเหมืองข้อมูลขั้นต่อไปได้

การพิจารณาความน่าสนใจ (interestingness) ของความรู้ที่ค้นพบโดยการทำเหมืองข้อมูล เป็นอีกประเด็นปัญหาสำคัญ ที่ได้รับความสนใจจากนักวิจัยจำนวนมากตั้งแต่ยุคเริ่มต้นของงานวิจัยด้านการทำเหมืองข้อมูล (Piatetsky-Shapiro and Matheus, 1994; Piatetsky-Shapiro et al., 1994) เนื่องจากได้เห็นถึงความสำคัญว่าความรู้ที่จะนำไปใช้ประโยชน์ได้ จะต้องเป็นความรู้ที่เป็นการค้นพบใหม่และน่าสนใจ โดย Piatetsky-Shapiro และ Matheus (1994) ได้พัฒนาระบบชื่อ KEFIR ขึ้นมาเพื่อใช้กับงานด้านการประกันสุขภาพ ระบบจะพิจารณาข้อมูลที่เบี่ยงเบนไปจากข้อมูลปกติ และตัดสินใจว่าข้อมูลที่เบี่ยงเบนนี้เป็นข้อมูลที่ควรค่าแก่การให้ความสนใจ การพิจารณาเกณฑ์ความน่าสนใจโดยระบบนี้จึงเป็นการพิจารณาในขอบเขตที่จำกัดมาก ในช่วงปีต่อมา Silberschatz และ Tuzhilin (1995; 1996) ได้ร่วมมือกันเสนอวิธีการพิจารณาความน่าสนใจของความรู้ที่ค้นพบโดยใช้ probabilistic belief เป็นแนวทางในการพิจารณาความน่าจะเป็นว่าความรู้ใดที่ผู้ใช้จะเห็นว่าน่าสนใจ

ปัจจัยที่นักวิจัยใช้พิจารณาประเด็นความน่าสนใจมีได้หลายปัจจัย ได้แก่ จำนวนข้อมูลที่ความรู้นั้นสามารถอธิบายได้ หรือที่เรียกว่า ความครอบคลุม (coverage) ความถูกต้องของความรู้ (confidence) ความไม่แปรผันของความรู้ (strength) ความมีนัยสำคัญของความรู้ (significance) ความเข้าใจง่ายของความรู้ (simplicity) ความใหม่ของการค้นพบ (unexpectedness) และการเป็นความรู้ที่นำไปสู่การปฏิบัติได้ (actionability) ปัจจัยเหล่านี้ถูกหยิบยกขึ้นมาพิจารณาโดยนักวิจัยหลายคน (Major and Mangano 1993; Piatetsky-Shapiro and Matheus 1994; Quinlan 1992)

ปัจจัยในด้าน coverage, confidence, strength, significance, simplicity เป็นปัจจัยที่สามารถสร้างเทคนิคขึ้นมาจัดการได้โดยไม่ต้องใช้ข้อมูลหรือคำแนะนำประกอบจากผู้ใช้ หรือจากความรู้อื่นที่เกี่ยวข้องกัน (domain knowledge) แต่ปัจจัยด้าน unexpectedness และ actionability เป็นปัจจัยที่พิจารณาคัดสินได้ยาก จึงจำเป็นต้องใช้ฮิวริสติก (heuristics) ช่วยในการพิจารณา

ในส่วนของจัดการกับปริมาณความรู้ ที่เป็นผลผลิตของกระบวนการทำเหมืองข้อมูล และมักจะได้รับความรู้ที่ไม่เป็นประโยชน์ปะปนออกมาเป็นจำนวนมาก เป็นปัญหาที่มีผลกระทบอย่างมากกับงานทำเหมืองข้อมูลโดยเฉพาะงานประเภท association (Bayardo et al., 1999; Brin et al., 1991; Frawley et al., 1991; Klemettinen et al., 1994; Suzuki 1997) นักวิจัยส่วนใหญ่จะใช้วิธีจัดการระหว่างกระบวนการทำเหมืองข้อมูล เช่น ใช้วิธีตัดกิ่ง (pruning) เส้นทางค้นหาที่คาดว่าจะนำไปสู่ความรู้ที่ไม่เกิดประโยชน์ นักวิจัยในกลุ่มนี้ได้แก่ Quinlan (1992), Breiman และคณะ (1984), Clark และ Matwin (1993) นักวิจัยในอีกกลุ่มจะใช้วิธีนำความรู้อื่น (domain knowledge) มาใช้ช่วยในระหว่างการค้นหาความรู้ โดยใช้สมมุติฐานว่าความรู้อื่นๆที่เกี่ยวข้องจะช่วยในการ

ตัดสินใจจัดตั้งผลลัพธ์ของการทำเหมืองข้อมูลที่จะไม่ก่อประโยชน์กับงาน นักวิจัยในกลุ่มนี้ได้แก่ Ortega และ Fisher (1995), Clark และ Matwin (1993), Pazzani และ Kibler (1992)

ถึงแม้แนวทางของเทคนิค pruning และเทคนิคการใช้ domain knowledge จะช่วยลดปริมาณความรู้ที่ไม่ก่อประโยชน์ แต่วัตถุประสงค์ของนักวิจัยในทั้งสองกลุ่มนี้มุ่งเน้นไปที่การลดปริมาณความรู้เพื่อเพิ่มประสิทธิภาพของกระบวนการค้นหาความรู้ (mining process) มากกว่าจะเพื่อเอื้อประโยชน์ให้กับผู้ใช้ ให้สามารถใช้ประโยชน์จากผลลัพธ์ของการทำเหมืองข้อมูลได้อย่างเต็มประสิทธิภาพและเกิดประโยชน์ในเวลาอันรวดเร็ว

ในระยะหลังของการวิจัยในสาขาการทำเหมืองข้อมูล นักวิจัยเริ่มตระหนักมากขึ้นว่ากระบวนการจัดการกับผลลัพธ์ (หรือ ความรู้) ที่ได้หลังจากการทำเหมืองข้อมูลเป็นสิ่งจำเป็นเพื่อให้เกิดการใช้ประโยชน์จากการทำเหมืองข้อมูลเป็นอัตโนมัติและทันต่อความต้องการมากขึ้น ดังจะเห็นได้จากงานวิจัยของ Wang และคณะ (1998), Adomavicius และ Tuzhilin (2001)

ซอฟต์แวร์ระบบเหมืองข้อมูลที่มีอยู่ในปัจจุบัน ส่วนใหญ่จะเน้นการพัฒนาเฉพาะส่วนค้นหาความรู้ (mining algorithms) โดยให้ขั้นตอนการเตรียมข้อมูลและการคัดเลือกข้อมูล (preprocessing) เป็นความรับผิดชอบของผู้ใช้ ซึ่งผู้ใช้ส่วนใหญ่จะต้องใช้เวลามากกว่า 80% ของการทำเหมืองข้อมูลในการเตรียมและคัดเลือกข้อมูล นอกจากนี้ซอฟต์แวร์ระบบเหมืองข้อมูลเกือบทั้งหมดจะแสดงผลลัพธ์ในปริมาณที่มากเกินไป ทำให้ผู้ใช้ต้องเป็นผู้เชี่ยวชาญด้านการวิเคราะห์ข้อมูลจึงจะสามารถแปลผล และคัดเลือกผลลัพธ์เพื่อนำไปใช้ประโยชน์ในงานบริหารจัดการหรือการตัดสินใจใดๆต่อไปได้

คณะผู้วิจัยจึงมีแนวคิดที่จะริเริ่มโครงการขนาดใหญ่ขึ้น ด้วยวัตถุประสงค์เพื่อนำไปสู่การพัฒนาระบบเหมืองข้อมูลที่มีขีดความสามารถสูงในชื่อของระบบเอสยูทีไมเนอร์ (SUT Miner) ที่มีความสมบูรณ์ครบถ้วน ในแง่ของการเตรียมฟังก์ชันที่จะอำนวยความสะดวกให้กับผู้ใช้ในการนำเข้าข้อมูล การเตรียมและคัดเลือกข้อมูลอย่างมีประสิทธิภาพ (ส่วน Preprocessing ในรูปที่ 1.2) มีฟังก์ชันการทำเหมืองข้อมูลแบบจัดกลุ่ม (clustering) และมีฟังก์ชันการทำเหมืองข้อมูลแบบจำแนก (classification) ที่สามารถนำไปใช้ประโยชน์ได้ในงานเพื่อการทำนาย (prediction) นอกจากนี้ยังได้จัดเตรียมฟังก์ชันกลั่นกรองความรู้ที่ได้จาก mining algorithms ให้ได้ผลลัพธ์ที่ตรงกับความต้องการของผู้ใช้มากที่สุดและง่ายต่อการแปลผลเพื่อนำไปสู่การใช้ประโยชน์จริง (ส่วน Postprocessing ในรูปที่ 1.2)

ระบบเหมืองข้อมูล SUT Miner ได้รับการพัฒนาขึ้นในลักษณะของโอเพนซอร์ส หรือเป็นซอฟต์แวร์ที่เปิดเผยซอร์สโค้ดและเป็นฟรีแวร์ เพื่อให้ผู้ใช้โดยทั่วไปสามารถนำไปใช้งาน หรือปรับปรุงซอร์สโค้ดให้ตรงกับความต้องการได้โดยไม่มีค่าใช้จ่าย ภาษาที่ใช้ในการพัฒนาโปรแกรมคือภาษาโปรล็อก ในรูปแบบของ SWI Prolog ที่เป็นฟรีแวร์เช่นเดียวกัน (ดาวน์โหลดได้จาก

www.swi-prolog.org) ผลสำเร็จของชุดโครงการนี้จึงคาดว่าจะช่วยให้ประเทศชาติประหยัดค่าใช้จ่ายในด้านการจัดซื้อซอฟต์แวร์ระบบเหมืองข้อมูลได้อย่างมาก และยังเป็นการพัฒนาเทคโนโลยีการทำเหมืองข้อมูลให้เอื้อประโยชน์ต่อระบบสนับสนุนการตัดสินใจ (decision support system) ได้มากขึ้น เนื่องจากผู้บริหาร โดยทั่วไปที่ไม่ใช่ นักวิเคราะห์ข้อมูลและนักคอมพิวเตอร์ จะสามารถใช้งานระบบเหมืองข้อมูลนี้ได้ง่ายขึ้นกว่าซอฟต์แวร์ที่มีอยู่ในปัจจุบัน

ผลสำเร็จของการพัฒนาส่วนกลั่นกรองความรู้ (postprocessing) จะเป็นนวัตกรรมใหม่ของเทคโนโลยีการทำเหมืองข้อมูล ให้ก้าวไปสู่การสร้างระบบผู้เชี่ยวชาญ (expert system) อย่างเป็นกระบวนการอัตโนมัติมากขึ้น เนื่องจากความรู้ที่เป็นผลลัพธ์จากการทำเหมืองข้อมูลจะได้รับการประเมินและกลั่นกรองความรู้ที่เป็นประโยชน์และแม่นยำตรง ส่งไปเก็บไว้ในฐานความรู้ (knowledge base) ของระบบผู้เชี่ยวชาญ ซึ่งจะช่วยให้งานของวิศวกรความรู้ (knowledge engineer) ง่ายขึ้นและรวดเร็วมากขึ้น

การดำเนินงานในชุดโครงการวิจัยนี้ นอกจากจะเป็นการสร้างนวัตกรรมด้านระบบเหมืองข้อมูล และช่วยพัฒนาความสามารถและทักษะของนักวิจัยแล้ว ยังเป็นการช่วยพัฒนาผู้ช่วยวิจัยให้เป็นบุคลากรที่มีความรู้ ความเชี่ยวชาญทางด้านเทคโนโลยีการทำเหมืองข้อมูลและการวิเคราะห์ข้อมูลอัตโนมัติ ซึ่งจะเป็นกำลังสำคัญในการพัฒนาประเทศชาติได้อย่างยั่งยืนต่อไป

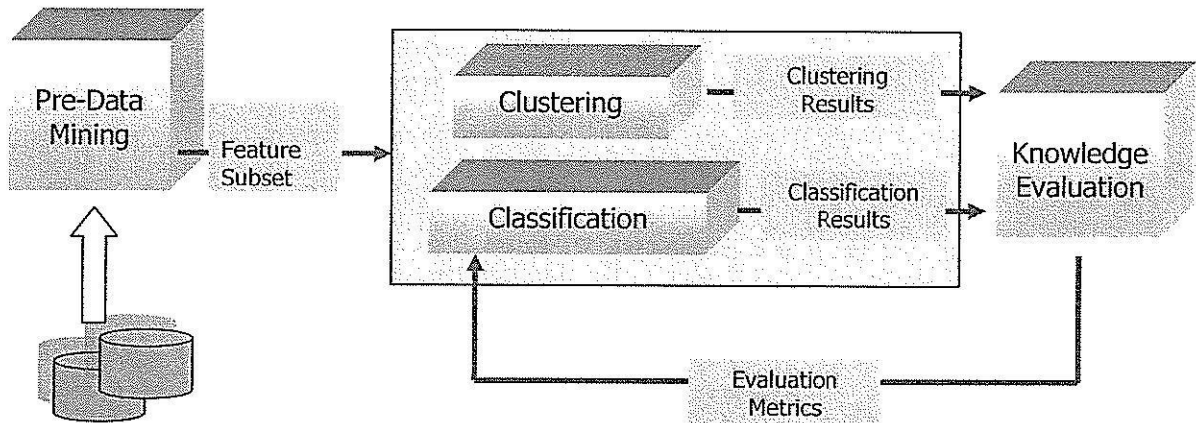
1.2 วัตถุประสงค์หลักของแผนงานวิจัย

- เพื่อพัฒนาระบบเหมืองข้อมูลที่มีความสมบูรณ์ครบถ้วนในลักษณะ integrated system ที่มีทั้งส่วน preprocessing ส่วน mining และส่วน postprocessing
- เพื่อออกแบบและพัฒนากระบวนการเตรียมข้อมูล และคัดเลือกข้อมูลให้เป็นวิธีการอัตโนมัติมากขึ้น
- เพื่อพัฒนา mining algorithms ประเภท clustering และ classification
- เพื่อพัฒนาแนวทางใหม่ในการคัดเลือกและกลั่นกรองความรู้ที่เป็นผลลัพธ์เบื้องต้นจากการทำเหมืองข้อมูล
- เพื่อสร้างเทคโนโลยีขึ้นใช้เองโดยไม่ต้องซื้อจากต่างประเทศ และเพิ่มความสามารถและทักษะให้กับนักวิจัยไทยให้สามารถแข่งขันได้ในระดับโลก

1.3 กรอบของแผนงานวิจัย

แผนงานวิจัยที่เสนอขึ้นนี้จะเน้นการพัฒนา ระบบเหมืองข้อมูลให้เป็นลักษณะระบบรวม (integrated system) กล่าวคือจะเป็นระบบที่ครอบคลุมส่วนประกอบหลักครบทั้งสามส่วน คือ ส่วน

การเตรียมข้อมูล ซึ่งเป็นขั้นตอนก่อนการทำเหมืองข้อมูล(pre-data mining) ส่วนการทำเหมืองข้อมูล(mining) ที่ทำหน้าที่ค้นหาความรู้ที่ซ่อนอยู่ในข้อมูลโดยจะมีวิธีการค้นหาความรู้ในสองรูปแบบ คือ แบบจัดกลุ่มข้อมูล (clustering) และแบบจำแนกประเภทข้อมูล (classification) ส่วนกลั่นกรองความรู้ (knowledge evaluation) จะเป็นขั้นตอนหลังการทำเหมืองข้อมูล (post-data mining) โดยแสดงส่วนประกอบทั้งหมดนี้เป็นแผนภาพได้ดังรูปที่ 1.3



รูปที่ 1.3 แผนภาพรวมของระบบเหมืองข้อมูล SUT Miner

ในการพัฒนาระบบเหมืองข้อมูลให้มีส่วนประกอบครบถ้วนสมบูรณ์ดังแผนภาพ จำเป็นต้องประกอบด้วยการพัฒนาในส่วนย่อย คือ

- ส่วนที่ 1 Pre-data mining ทำหน้าที่คัดเลือกข้อมูลและปรับปรุงคุณภาพข้อมูล (โครงการวิจัยที่ 1)
- ส่วนที่ 2 Clustering data mining ทำหน้าที่จัดกลุ่มข้อมูลโดยอัตโนมัติ (โครงการวิจัยที่ 2)
- ส่วนที่ 3 Classification data mining ทำหน้าที่ค้นหารูปแบบที่สามารถใช้จำแนกข้อมูล (โครงการวิจัยที่ 3)
- ส่วนที่ 4 Post-data mining ทำหน้าที่ประเมินผลลัพธ์ที่ได้และคัดเลือกแสดงผลเฉพาะผลลัพธ์ที่คาดว่าจะประโยชน์และน่าสนใจสำหรับผู้ใช้ (โครงการวิจัยที่ 4)

ซอฟต์แวร์ระบบเหมืองข้อมูลที่ตีพิมพ์มีอยู่ในปัจจุบันจะมีส่วนที่ 1, 2, และ 3 การพัฒนาเพิ่มเติมส่วนที่ 4 ในงานวิจัยนี้จะเป็นแนวทางใหม่ของการสร้างระบบเหมืองข้อมูล ที่คาดหมายว่าจะเป็นประโยชน์สำหรับผู้ใช้มากขึ้น