



รายงานวิจัยฉบับสมบูรณ์

โครงการ

การค้นพบโดยประมาณของรูปแบบที่ปรากฏบ่อยในข้อมูลสตรีม
Approximate Frequent Pattern Discovery over Data Stream

โดย รองศาสตราจารย์ ดร.กิตติศักดิ์ เกิดประสพ

20 สิงหาคม 2553

600251713

ห้องสมุดงานวิจัย สำนักงานคณะกรรมการการวิจัยแห่งชาติ



246593

สัญญา

รายงานวิจัยฉบับสมบูรณ์

โครงการ

การค้นพบโดยประมาณของรูปแบบที่ปรากฏบ่อยในข้อมูลสตรีม
Approximate Frequent Pattern Discovery over Data Stream



ผู้วิจัย

รองศาสตราจารย์ ดร.กิตติศักดิ์ เกิดประสพ

สาขาวิชาวิศวกรรมคอมพิวเตอร์

สำนักวิชาวิศวกรรมศาสตร์

มหาวิทยาลัยเทคโนโลยีสุรนารี

สนับสนุนโดยสำนักงานคณะกรรมการการอุดมศึกษา

และสำนักงานกองทุนสนับสนุนการวิจัย

(ความเห็นในรายงานนี้เป็นของผู้วิจัย สกอ. และ สกว. ไม่จำเป็นต้องเห็นด้วยเสมอไป)

กิตติกรรมประกาศ

ผู้วิจัยขอขอบคุณสำนักงานคณะกรรมการการอุดมศึกษา (สกอ.) และสำนักงานกองทุนสนับสนุนการวิจัย (สกว.) ที่ได้เล็งเห็นถึงความสำคัญของการเสริมสร้างนักวิจัยและพัฒนางานวิจัยของชาติ โดยทั้งสองหน่วยงานได้ร่วมมือกันจัดสรรทุนวิจัยให้กับนักวิจัยที่ทำงานอยู่ในมหาวิทยาลัยต่างๆ ทั่วประเทศอย่างต่อเนื่อง ความสำเร็จของงานวิจัยนี้จะเกิดขึ้นไม่ได้หากไม่ได้รับการสนับสนุนจากผู้บริหารของหน่วยงานต้นสังกัดได้แก่รองศาสตราจารย์ นาวาอากาศเอก ดร.วรพจน์ ขำพิศ คณบดีสำนักวิชาวิศวกรรมศาสตร์ และศาสตราจารย์ ดร.ประสพ สืบคำ อธิการบดีมหาวิทยาลัยเทคโนโลยีสุรนารี ที่ตระหนักถึงความสำคัญของการสร้างนักวิจัย และสนับสนุนคณาจารย์ในมหาวิทยาลัยให้ได้ทำงานวิจัยอย่างต่อเนื่อง โดยได้เตรียมสภาพแวดล้อมในการทำงานทั้งสถานที่และอุปกรณ์การวิจัยที่เหมาะสมสำหรับนักวิจัย รวมถึงการสนับสนุนให้คณาจารย์ได้มีโอกาสดำเนินงานไปเผยแพร่ผลงานวิจัยต่อที่ประชุมวิชาการทั้งในและต่างประเทศ ผู้วิจัยขอขอบคุณทุกท่านไว้ ณ โอกาสนี้

บทคัดย่อ

รหัสโครงการ: RMU5080026

ชื่อโครงการ: การค้นพบโดยประมาณของรูปแบบที่ปรากฏบ่อยในข้อมูลสตรีม

ชื่อนักวิจัย: รศ.ดร.กิตติศักดิ์ เกิดประสพ

สาขาวิชาวิศวกรรมคอมพิวเตอร์ มหาวิทยาลัยเทคโนโลยีสุรนารี

E-mail address: kerdpras@sut.ac.th, KittisakThailand@gmail.com

ระยะเวลาโครงการ: 25 มิถุนายน 2550 – 24 มิถุนายน 2553

246593

การค้นพบรูปแบบที่ปรากฏบ่อย เป็นปฏิบัติการที่สำคัญในงานวิเคราะห์ความสัมพันธ์ กระบวนการค้นพบนี้เป็นการคัดแยกโดยอัตโนมัติ เพื่อค้นหารูปแบบที่น่าสนใจและความสัมพันธ์ที่เกิดขึ้นร่วมกันในฐานข้อมูลขนาดใหญ่ รูปแบบที่น่าสนใจเหล่านี้สามารถนำไปใช้เพื่อการสร้างกฎความสัมพันธ์สำหรับสนับสนุนการตัดสินใจในงานด้านต่างๆ เช่น การพยากรณ์ด้านการเงิน การวินิจฉัยทางการแพทย์ งานวิจัยปัจจุบันด้านการทำเหมืองข้อมูลเพื่อค้นหาความสัมพันธ์ จะมุ่งความสนใจไปที่การพัฒนาวิธีการที่มีประสิทธิภาพสูงเพื่อค้นหาวัตถุหรือไอเท็มที่ปรากฏร่วมกัน ทั้งนี้เนื่องจากถ้าในฐานข้อมูลประกอบด้วยไอเท็มที่แตกต่างกันจำนวน m รายการ ไอเท็มเหล่านี้สามารถทำให้เกิดรูปแบบร่วมได้มากถึง 2^m รูปแบบ ซึ่งการค้นหาแบบร่วมจำนวนมากเช่นนี้จะเป็งานประมวลผลที่ใช้เวลานานมาก การค้นหาที่ยังยากขึ้นถ้าข้อมูลเป็นลักษณะสตรีม เนื่องจากลักษณะของสตรีมจะผลิตข้อมูลอย่างต่อเนื่องในปริมาณมาก ทำให้การทำงานกับข้อมูลต้องกระทำในรอบเดียวเพื่อให้ได้ผลการวิเคราะห์ความสัมพันธ์ที่ทันต่อการใช้งาน จากการศึกษาถึงลักษณะของปัญหาและข้อจำกัดต่างๆเหล่านี้ ผู้วิจัยจึงได้พัฒนาวิธีการโดยประมาณเพื่อค้นหาแบบที่ปรากฏบ่อยในข้อมูลสตรีม วิธีการโดยประมาณจะถูกนำมาใช้กับข้อมูลสตรีมก่อนที่จะนำข้อมูลตัวแทนไปประมวลผลต่อเพื่อค้นหาแบบที่ปรากฏบ่อย วิธีการที่เสนอขึ้นนี้ได้รับการพัฒนาเป็นโปรแกรมต้นแบบด้วยภาษาเชิงฟังก์ชัน ผลการทดสอบโปรแกรมกับข้อมูลจริงได้ผลลัพธ์ที่น่าพอใจในด้านประสิทธิภาพและความถูกต้องของรูปแบบที่เป็นผลลัพธ์ของโปรแกรม แนวทางการพัฒนาต่อยอดของงานวิจัยนี้ จะเป็นการใช้วิธีการประมวลผลโปรแกรมแบบคู่ขนานที่คาดว่าจะช่วยให้สามารถทำงานกับข้อมูลขนาดใหญ่มากได้

คำหลัก: การค้นพบรูปแบบที่ปรากฏบ่อย, ข้อมูลสตรีม, วิธีการโดยประมาณ

Abstract

Project Code: RMU5080026

Project Title: Approximate frequent pattern discovery over data stream

Investigator: Dr. Kittisak Kerdprasop, Associate Professor
School of Computer Engineering, Suranaree University of Technology

E-mail Address: kerdpras@sut.ac.th, KittisakThailand@gmail.com

Project Period: 25 June 2007 – 24 June 2010

246593

Frequent pattern discovery is an essential operation for association analysis. The discovery process concerns an automatic extraction of interesting patterns and correlations from a large database. These patterns can reveal implicit relationships among set of objects (or items) that lead to the generation of association rules to be used for decision support, financial forecast, medical diagnosis and many other applications. Current studies in association rule mining concentrate on how to effectively find all objects frequently co-occurring. Given m objects, there are as much as 2^m frequent patterns to consider. Frequent pattern discovery is thus a computationally expensive problem. It is even harder over data stream because a continuously generated nature of stream does not allow a revisit on each data element. Furthermore, pattern discovery process must be fast to produce timely results. Based on these requirements, we devise an approximate approach to tackle the problem of discovering frequent patterns over continuous stream. Our approximation algorithm is intended to be applied to process a stream prior to the pattern discovery process. We propose a stochastic method to get a good guess of the stream characteristics, and then draw a set of representatives from the incoming stream. These representatives are subsequently used in the process of frequent pattern mining. Our design had been implemented with the functional programming paradigm and the experimental results confirm the efficiency and reliability of our method. For a massive database, parallel method is a solution for the scalability problem. That is the main direction of our future research.

Keywords: Frequent pattern discovery, Data stream, Approximation method

บทสรุปสำหรับผู้บริหาร

Executive Summary

รหัสโครงการ: RMU5080026

ชื่อโครงการ: การค้นพบโดยประมาณของรูปแบบที่ปรากฏบ่อยในข้อมูลสตรีม
Approximate frequent pattern discovery over data stream

ชื่อนักวิจัย: รศ.ดร.กิตติศักดิ์ เกิดประสพ
สาขาวิชาวิศวกรรมคอมพิวเตอร์ มหาวิทยาลัยเทคโนโลยีสุรนารี

E-mail address: kerdpras@sut.ac.th, KittisakThailand@gmail.com

ระยะเวลาโครงการ: 25 มิถุนายน 2550 – 24 มิถุนายน 2553

ปัญหาที่ทำการวิจัย และความสำคัญของปัญหา

การค้นพบความสัมพันธ์หรือความเชื่อมโยงของสินค้าที่ถูกซื้อพร้อมกันบ่อยเช่น ลูกค้ายี่ห้อเบียร์ มักจะซื้อผ้าอ้อมเด็กด้วย เรียกว่า การทำเหมืองเพื่อค้นหาความสัมพันธ์ (association mining) ซึ่งเป็นวิธีการทำเหมืองข้อมูลประเภทหนึ่ง การค้นพบความสัมพันธ์อาศัยขั้นตอนพื้นฐานที่สำคัญคือ การค้นพบรูปแบบที่ปรากฏบ่อย (frequent pattern discovery) รูปแบบนี้ได้หลากหลายประเภท เช่น รูปแบบการซื้อสินค้าของลูกค้าส่วนใหญ่ในห้างสรรพสินค้า, รูปแบบการเกิดแผ่นดินไหวในรอบสิบปีที่ผ่านมาของภูมิภาคเอเชียตะวันออกเฉียงใต้, รูปแบบการค้นหาข้อมูลจากเว็บเพจต่างๆของผู้ใช้อินเทอร์เน็ต, รูปแบบการจัดโครงสร้างโมเลกุลในสารประกอบโปรตีน เป็นต้น

นับตั้งแต่ปีค.ศ.1993 ที่ได้มีการเริ่มเสนอแนวคิดและเทคนิคของการทำเหมืองข้อมูลประเภทการค้นพบรูปแบบที่ปรากฏบ่อย งานวิจัยด้านนี้ก็ได้รับความนิยมอย่างกว้างขวางเนื่องจากใช้ประโยชน์ได้กับการวิเคราะห์ข้อมูลในหลากหลายสาขา มีการพัฒนาเทคนิคต่างๆของการค้นหารูปแบบที่ปรากฏบ่อยให้มีประสิทธิภาพสูงขึ้น ในระยะหลังความก้าวหน้าของเทคโนโลยีอินเทอร์เน็ตทำให้รูปแบบของข้อมูลเปลี่ยนแปลงไปจากข้อมูลที่มีขนาดใหญ่แต่คงตัว (static) กลายเป็นข้อมูลที่มีลักษณะพลวัต (dynamic) เปลี่ยนแปลงได้ทั้งการเพิ่มปริมาณอย่างต่อเนื่องและการเปลี่ยนรูปแบบการกระจายของข้อมูล ข้อมูลลักษณะเช่นนี้เรียกว่า ข้อมูลสตรีม (data stream)

ในช่วงเวลาห้าปีที่ผ่านมา นักวิจัยจำนวนหนึ่งในสาขาการทำเหมืองข้อมูล เริ่มให้ความสนใจกับข้อมูลสตรีม พยายามปรับปรุงเทคนิคการค้นหารูปแบบที่ปรากฏบ่อยที่ใช้งานได้ดีกับข้อมูลคงตัวให้ทำงานกับข้อมูลสตรีมได้ แต่เทคนิคเหล่านั้นยังมีข้อจำกัดและใช้ได้กับเพียงบางแอปพลิเคชัน ผู้วิจัยจึงได้เสนอโครงการวิจัยนี้ขึ้นเพื่อพัฒนาเทคนิคการค้นหารูปแบบที่ปรากฏบ่อยจากข้อมูลสตรีม ที่ใช้งานได้กับหลากหลายแอปพลิเคชัน เช่น web-page click stream, DNA stream, market basket data stream และเพื่อให้ได้ผลลัพธ์ของการค้นหารูปแบบที่ปรากฏบ่อยในเวลาที่รวดเร็วเพื่อทันต่อความต้องการใช้งาน

การประมวลผลข้อมูลสตรีมจะใช้การประมาณค่าด้วยกระบวนการ stochastic ในขั้นตอนของการพัฒนาโปรแกรมต้นแบบจะใช้แนวทางการโปรแกรมเชิงฟังก์ชัน ที่มีโครงสร้างของภาษานับสนุนการตรวจจับแพทเทิร์น การพัฒนาเทคนิคการค้นพบรูปแบบที่ปรากฏบ่อยในข้อมูลสตรีมโดยใช้ stochastic process วิธีการค้นพบโดยประมาณและการทำ rapid prototyping ด้วยการใช้การโปรแกรมเชิงฟังก์ชัน รวมถึงแนวทางการทดสอบประสิทธิภาพของโปรแกรม จะเป็นองค์ความรู้ใหม่ที่ช่วยพัฒนาความก้าวหน้าของงานวิจัยในสาขาการทำเหมืองข้อมูลประเภทการค้นพบรูปแบบที่ปรากฏบ่อย และการค้นพบความสัมพันธ์

วัตถุประสงค์

- เพื่อพัฒนาเทคนิคและอัลกอริทึมในการค้นพบโดยประมาณของรูปแบบที่ปรากฏบ่อยในข้อมูลสตรีม
- พัฒนาการอบงาน stochastic process ที่เหมาะสมกับการทำ association mining กับข้อมูลสตรีม
- ออกแบบและพัฒนาโปรแกรมต้นแบบในการทำ approximate association mining ด้วยเทคนิคการโปรแกรมเชิงฟังก์ชัน

ผลสำเร็จของโครงการ

1. ผลงานตีพิมพ์ในหนังสือวิชาการ (Book chapters)

- 1) K. Kerdprasop and N. Kerdprasop (2009). Knowledge mining with a higher-order logic approach. In K. Nakamatsu, G. Phillips-Wrens, L.C. Jain, R.J. Howlett (Eds.), *New Advances in Intelligent Decision Technologies*, pp. 151-159, Springer. (ISBN: 978-3-642-00908-2, DOI: 10.1007/978-3-642-00909-9)
- 2) N. Kerdprasop, S. Pilabutr, and K. Kerdprasop (2009). Improving medical database consistency with induced trigger rules. In K. Nakamatsu, G. Phillips-Wrens, L.C. Jain, R.J. Howlett (Eds.), *New Advances in Intelligent Decision Technologies*, pp. 265-274, Springer. (ISBN: 978-3-642-00908-2, DOI: 10.1007/978-3-642-00909-9)

2. ผลงานตีพิมพ์ในวารสารวิชาการนานาชาติ (Journal articles)

- 1) N. Pannurat, N Kerdprasop, and K. Kerdprasop (2010). Database reverse engineering based on association rule mining. *International Journal of Computer Science Issues (IJCSI)*, Volume 7, Issue 2, March 2010, pp. 10-15.
- 2) N. Kerdprasop and K. Kerdprasop (2009). Knowledge induction from medical databases with higher-order programming. *WSEAS Transactions on Information Science and Applications*, Issue 10, Volume 6, October 2009, pp. 1719-1728.
- 3) K. Kerdprasop and N. Kerdprasop (2008). Approximate frequent pattern discovery over data stream. *International Journal of Computer Science and Engineering (IJCSE)*, Volume 2, Number 1, pp. 28-32.
- 4) K. Kerdprasop and N. Kerdprasop (2007). On pattern-based programming towards the discovery of frequent patterns. *International Journal of Computer Science (IJCS)*, Volume 2, Number 4, pp. 268-273.

- 5) N. Kerdprasop and K. Kerdprasop (2007). Mining frequent patterns with functional programming. *International Journal of Computer and Information Science and Engineering (IJCISE)*, Volume 1, Number 2, pp. 66-71.
- 6) N. Kerdprasop and K. Kerdprasop (2007). Moving data mining tools toward a business intelligence system. *International Journal of Intelligent Technology (IJIT)*, Volume 2, Number 2, pp. 99-104.

Cited by:

- Z. Zhu, J. Gu, L. Zhang, W. Song, and R. Gao. (2009). Research on domain-driven actionable knowledge discovery. *Proc. 20th Int. Conf. on Cutting-Edge Research Topics on Multiple Criteria Decision Making (MCDM 2009)*, 21-26 June 2009, China, pp.176-183.
- M. Al-Noukari and W. Al-Hussan. (2007). Using data mining techniques for predicting future car market demand: dex case study. *Proc. 3rd Int. Conf. on Information and Communication Technology (ICTTA 2008)*, 7-11 April 2008, pp.1-5.

3. การเสนอผลงานในที่ประชุมวิชาการและตีพิมพ์ผลงานใน Conference Proceedings

- 1) K. Kerdprasop and N. Kerdprasop (2009). Automated induction of frequent patterns with knowledge-based software engineering. *Proceedings of the Joint Conference of ASCM 2009 (Asian Symposium on Computer Mathematics) and MACIS 2009 (Mathematical Aspects of Computer and Information Sciences)*, 14-17 December 2009, Fukuoka, Japan, pp. 431-434.
- 2) N. Kerdprasop and K. Kerdprasop (2008). A declarative programming paradigm and the development of knowledge mining agents. *Proceedings of IADIS Multi Conference on Computer Science and Information Systems*, Amsterdam, Netherlands, 22-27 July 2008, pp. 45-52.

4. การนำผลงานวิจัยไปใช้ประโยชน์เชิงวิชาการ

- 1) ใช้ประกอบการเรียนการสอนในรายวิชา 423681 Selected Topics in Computer Engineering I ซึ่งเป็นรายวิชาเลือกในหลักสูตรวิศวกรรมศาสตรมหาบัณฑิต สาขาวิชาวิศวกรรมคอมพิวเตอร์
- 2) สร้างนักวิจัยใหม่ด้วยการนำส่วนหนึ่งของงานวิจัยไปเป็นหัวข้อวิทยานิพนธ์ “การปรับรูปแบบบรรทัดฐานในฐานข้อมูลเชิงสัมพันธ์ด้วยเทคนิคการวิเคราะห์ความสัมพันธ์” (Normalization in relational database with association analysis technique) โดยนาย ณัฐพล พันนุรัตน์ นักศึกษาในหลักสูตรวิศวกรรมศาสตรมหาบัณฑิต สาขาวิชาวิศวกรรมคอมพิวเตอร์

สารบัญ

	หน้า
กิตติกรรมประกาศ	ii
บทคัดย่อภาษาไทย	iii
บทคัดย่อภาษาอังกฤษ	iv
บทสรุปสำหรับผู้บริหาร (Executive summary)	v
สารบัญ	viii
สารบัญภาพ	ix
สารบัญตาราง	x
บทที่ 1 บทนำ	1
1.1 ความสำคัญและที่มาของปัญหาที่ทำการวิจัย	1
1.2 วัตถุประสงค์ของโครงการ	2
1.3 ผลงานวิจัยที่เกี่ยวข้อง	2
1.4 แนวทางการวิจัย	7
บทที่ 2 การออกแบบและพัฒนาโปรแกรม	8
2.1 การออกแบบวิธีการค้นพบโดยประมาณของรูปแบบที่ปรากฏบ่อย	8
2.2 การค้นพบรูปแบบที่ปรากฏบ่อยด้วยภาษาสาสเกล	12
2.3 การพัฒนาโปรแกรมต้นแบบด้วยภาษาเออแลง	13
บทที่ 3 การทดสอบโปรแกรมและผลการทดสอบ	19
3.1 วิธีการทดสอบโปรแกรม	19
3.2 ข้อมูลที่ใช้และผลการทดสอบ	20
บทที่ 4 บทสรุปและวิจารณ์	26
4.1 สรุปผลการวิจัย	26
4.2 แนวทางการพัฒนาในอนาคต	28
เอกสารอ้างอิง	29
ผลผลิตจากโครงการวิจัยที่ได้รับทุนจากสกอ. และ สกว.	32
ภาคผนวก	
ก บทความดีพิมพ์จากโครงการวิจัย	34
ข บทความสำหรับการเผยแพร่	109

สารบัญภาพ

	หน้า
รูปที่ 1.1 . โครงข่ายของไอเท็มเซตทั้งหมดที่ต้องพิจารณาเพื่อค้นหาไอเท็มเซตที่ปรากฏบ่อย	4
รูปที่ 1.2 โครงข่ายของไอเท็มเซตที่ขนาดเล็กลงจากการไม่ต้องพิจารณาไอเท็ม E	5
รูปที่ 2.1 กรอบของงานออกแบบวิธีการค้นพบรูปแบบที่ปรากฏบ่อยในข้อมูลสตรีม	8
รูปที่ 2.2 อัลกอริทึมการเลื่อนกรอบหน้าต่างเพื่อคำนวณความหนาแน่นของข้อมูล	9
รูปที่ 2.3 การกำหนดกรอบหน้าต่างและขอบเขตการนับจำนวนข้อมูลในแต่ละกรอบหน้าต่าง	10
รูปที่ 2.4 ข้อมูลสรุปของกรอบหน้าต่างที่ 16, 17, 29, 30 และ 35	10
รูปที่ 2.5 การแปลงข้อมูลกลับเป็นเวกเตอร์ที่แทนจุดข้อมูล	11
รูปที่ 2.6 อัลกอริทึมการสุ่มข้อมูลตามความหนาแน่นและการ maintain ข้อมูล	12
รูปที่ 2.7 ตัวอย่างโปรแกรมภาษาฮาสเกิลเพื่อคำนวณค่าไฟโบนาซชี	13
รูปที่ 2.8 โปรแกรมภาษาฮาสเกิลเพื่อการค้นพบรูปแบบที่ปรากฏบ่อย	13
รูปที่ 2.9 การเขียนฟังก์ชันแฟลคทอเรียลเปรียบเทียบระหว่างภาษาฮาสเกิลและภาษาเออแลง	14
รูปที่ 3.1 วิธีการทดสอบความสมบูรณ์และถูกต้องของการค้นพบรูปแบบที่ปรากฏบ่อยในข้อมูลสตรีม	19
รูปที่ 3.2 ตัวอย่างข้อมูลรหัสพันธุกรรม	20
รูปที่ 3.3 โครงสร้างของสายรหัสพันธุกรรมที่ประกอบด้วยส่วน exon และ intron	20
รูปที่ 3.4 ตัวอย่างการใช้โปรแกรมค้นพบรูปแบบที่ปรากฏบ่อยในข้อมูลรหัสพันธุกรรม	21
รูปที่ 3.5 ตัวอย่างการใช้ฟังก์ชัน findPos แปลความหมายรูปแบบรหัสพันธุกรรม	22
รูปที่ 3.6 การใช้ฟังก์ชัน findSupOf เพื่อตรวจสอบรูปแบบที่ปรากฏบ่อยในชุดข้อมูลทดสอบ	23
รูปที่ 4.1 การค้นพบรูปแบบโดยประมาณจากข้อมูลสตรีม	28

สารบัญตาราง

	หน้า
ตารางที่ 1.1 ข้อมูลตัวอย่างแสดงรายการสินค้าที่ลูกค้าซื้อจากร้านค้า	3
ตารางที่ 3.1 เปรียบเทียบจำนวนแพทเทิร์นที่ค้นพบจากข้อมูลสตรีม	24
ตารางที่ 3.2 เปรียบเทียบสัดส่วนการค้นพบแพทเทิร์นจากข้อมูลทดสอบ	25