

## บทคัดย่อ

รหัสโครงการ: MRG5380130

ชื่อโครงการ: การจัดกลุ่มกระแสข้อมูลอนุกรรมเวลาอย่างมีความหมายและแม่นยำ

ผู้วิจัย: ผู้ช่วยศาสตราจารย์ โอดิรัตน์ รัตนาหทธรน และคณะ  
ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย

อีเมล: chotirat@gmail.com, chotirat.r@chula.ac.th

ระยะเวลาโครงการ: 2 ปี

### บทคัดย่อ:

การจัดกลุ่มลำดับย่อยสำหรับข้อมูลอนุกรรมเวลาแบบกระแส เป็นหนึ่งในปัญหาที่ท้าทายมากที่สุด ของการทำเหมืองข้อมูลอนุกรรมเวลา ตั้งแต่การจัดกลุ่มลำดับย่อยได้ถูกแสดงให้เห็นว่า การจัดกลุ่มจะให้ คำตอบที่เร็วความหมายในเชิงการทดสอบ และทฤษฎี การจัดกลุ่มลำดับย่อยของข้อมูลอนุกรรมเวลาที่ถูกใช้ ในหลายร้อยงานวิจัยนั้นจะให้คลื่นไชน์เป็นตัวแทนกลุ่มเสมอ ถ้าให้ข้อมูลอนุกรรมเวลาหนึ่ง ๆ การจัดกลุ่ม ลำดับย่อยของข้อมูลอนุกรรมเวลาควรค้นค่าตัวแทนกลุ่มที่เป็นลักษณะของทุกลำดับย่อยในข้อมูลอนุกรรม เวลา ส่วนสาเหตุที่ทำให้เกิดความเร็วความหมาย ถูกระบุไว้เป็นสองสาเหตุได้แก่ การใช้ระยะทางยุคคลิดเป็น ตัววัดระยะทางที่ไม่เหมาะสม และการใช้การเฉลี่ยค่าตามแม่พลิจูดเป็นฟังก์ชันการเฉลี่ยที่ไม่เหมาะสม เพื่อที่จะได้มาซึ่งคำตอบของการจัดกลุ่มที่มีความหมาย ในงานวิจัยนี้ได้เสนอการจัดกลุ่มลำดับย่อยของ ข้อมูลอนุกรรมเวลาตามรูป โดยใช้ระยะทางโดยนิยมใหม่儿ร์บปิงและการเฉลี่ยค่าตามรูปแพนระยะทางยุค คลิด และการเฉลี่ยค่าตามแม่พลิจูดตามลำดับ ดังนั้นการจัดกลุ่มลำดับย่อยของข้อมูลอนุกรรมเวลาตามรูป จะคืนผลลัพธ์ที่มีความหมายที่มากกว่าการจัดกลุ่มลำดับย่อยของข้อมูลอนุกรรมเวลาแบบเดิม แต่อย่างไรก็ ตาม การจัดกลุ่มลำดับย่อยของข้อมูลอนุกรรมเวลาตามรูปไม่สามารถประยุกต์ใช้กับข้อมูลแบบกระแสได้ โดยตรง เนื่องจากการจัดกลุ่มลำดับย่อยของข้อมูลอนุกรรมเวลาตามรูปใช้เวลาในการประมวลผลนาน โดย คำนวณลำดับย่อยที่ผ่านมาทั้งหมด เมื่อมีจุดข้อมูลใหม่เข้ามา งานวิจัยนี้จึงได้เสนอการจัดกลุ่มลำดับย่อย ของข้อมูลอนุกรรมเวลาแบบกระแสตามรูป ให้รองรับกรณีข้อมูลแบบกระแส โดยคำนวณบนชุดข้อมูล ขนาดเล็กของลำดับย่อยที่เก็บไว้ แทนที่จะคำนวณจากลำดับย่อยทั้งหมด ซึ่งชุดข้อมูลของลำดับย่อยที่เก็บ ไว้ถูกปรับปรุงสำหรับทุก ๆ จุดข้อมูล เพื่อรักษาจำนวนลำดับย่อยในชุดข้อมูล ไม่ให้เกินกว่าจำนวนมากสุด ที่อนุญาต ดังนั้นการจัดกลุ่มลำดับย่อยของข้อมูลอนุกรรมเวลาแบบกระแสตามรูป จึงเร็วกว่าการจัดกลุ่ม ลำดับย่อยของข้อมูลอนุกรรมเวลาตามรูปอย่างมาก

คำหลัก : การทำเหมืองข้อมูล, การจัดกลุ่มลำดับย่อย, อนุกรรมเวลา, ข้อมูลแบบกระแส

## **Abstract**

---

**Project Code : MRG5380130**

**Project Title : Meaningful and Accurate Subsequence Clustering for Time Series Data Stream**

**Investigator : Assistant Professor Chotirat Ratanamahatana, Ph.D.**

**Dept. of Computer Engineering, Faculty of Engineering, Chulalongkorn University**

**E-mail Address : chotirat@gmail.com, chotirat.r@chula.ac.th**

**Project Period : 2 years**

### **Abstract:**

Subsequence clustering for time series data streams is one of the most challenging issues of time series data mining since subsequence clustering has been proven both theoretically and empirically that it produces meaningless clustering results, where hundreds of research works that utilize STSC as a preprocessing step and a subroutine are all affected. Given a time series sequence, subsequence clustering should return cluster representatives which represent characteristics of all subsequences in time series. Therefore, if cluster representatives are always sine waves regardless of inputs, clustering results are meaningless since they do not reflect characteristics of the subsequences. The causes of meaninglessness are identified in twofold, i.e., inappropriate uses of Euclidean distance as a distance measure and Amplitude Averaging as an averaging function. To achieve meaningful clustering results, in this research, Shape-based Subsequence Time Series Clustering (2STSC) is proposed to use Dynamic Time Warping (DTW) distance measure and Shape-based Averaging function. Therefore, 2STSC returns more meaningful results than those from STSC. However, 2STSC cannot directly apply to data streams since 2STSC consumes large computational complexity by considering all previous subsequences for every new incoming data point. Shape-based Streaming Subsequence Time Series Clustering (3STSC) is then proposed to handle the streaming case by calculating a clustering result on a small set of stored subsequences instead of calculating from all previous subsequences. The small set of stored subsequences is updated for every new incoming data point to maintain the number of stored subsequences not to exceed the maximum allowance number. 3STSC, therefore, is much faster than 2STSC, while 3STSC returns small distortions of clustering results.

**Keywords : Data Mining, Subsequence Clustering, Time Series, Data Stream**