

การจัดกลุ่มบทความทางวิทยาศาสตร์ของยีสต์

Clustering of Scientific Abstracts of Yeast

จรรย์ยุทธ ไชยจรรวมวิช ธิติพล เทียมจันทร์

ภาควิชาวิทยาการคอมพิวเตอร์ คณะวิทยาศาสตร์ มหาวิทยาลัยเชียงใหม่

อ.เมือง จ. เชียงใหม่ 50202 โทรศัพท์ 053-943409 โทรสาร 053-943433

Email: jeerayut@science.cmu.ac.th compsci18@hotmail.com

ABSTRACT - This paper uses the covariance-based clustering proposed by the author in [1] to analyze the scientific abstract corpus of yeast given in KDD Cup 2002 (Task 2) competition to show the effectiveness of the algorithm even dealing with data of large size.

KEY WORDS - Text clustering, Spectral data clustering

บทคัดย่อ - บทความนี้ นำวิธีการจัดกลุ่มสมาชิกของเมตริกซ์จำนวนจริงแบบไม่สมมาตร โดยวิธี covariance-based clustering ที่เสนอใน [1] มาใช้ในการวิเคราะห์กลุ่มเอกสารบทความทางวิทยาศาสตร์ของยีสต์ ที่กำหนดให้ในการแข่งขันการสืบค้นความรู้ KDD Cup 2002 (Task2) เพื่อแสดงให้เห็นว่าวิธีการจัดกลุ่มดังกล่าวสามารถประยุกต์ใช้งานได้แม้ข้อมูลจะมีขนาดใหญ่

คำสำคัญ - การจัดกลุ่มเอกสารข้อความ การจัดกลุ่มข้อมูลด้วยวิธีสเปกตรัล

1. บทนำ

การทำเหมืองข้อมูลมีจุดประสงค์หลักคือ การวิเคราะห์ค้นหาคำรู้หรือรูปแบบความสัมพันธ์ที่แฝงอยู่ในข้อมูลที่สามารถนำไปใช้ให้เกิดประโยชน์ ในปัจจุบันมีเทคนิคที่ใช้วิเคราะห์ข้อมูลหลายเทคนิคด้วยกัน ได้แก่ เทคนิคการแบ่งกลุ่มข้อมูล (data clustering) เทคนิคการหาความสัมพันธ์ (association rule) เทคนิคการจัดกลุ่มและทำนายกลุ่มข้อมูล (data classification and prediction) [2] เป็นต้น อีกทั้งยังมีการนำเอาหลายเทคนิคไปประยุกต์ใช้ร่วมกันเพื่อให้เหมาะสมกับงานต่างๆ สำหรับเทคนิคการแบ่งกลุ่มหรือคลัสเตอร์ข้อมูลนั้น เป็นเทคนิคที่ได้รับความนิยมในการนำไปประยุกต์ใช้กับเทคนิคอื่น และเหมาะที่จะนำไปแบ่งกลุ่มข้อมูลในกรณีที่ไม่มีการระบุประเภทหรือจำนวนกลุ่มที่มีอยู่ของข้อมูลมาก่อน โดยหลักการของเทคนิคคลัสเตอร์ข้อมูลคือ การจัดข้อมูลให้เป็นสมาชิกของกลุ่มเดียวกันหากมีความคล้ายคลึงกัน และ ข้อมูลที่เป็นสมาชิกของกลุ่มหนึ่งจะไม่คล้ายคลึงกับข้อมูลที่เป็นสมาชิกของกลุ่มอื่น

เทคนิคคลัสเตอร์ข้อมูลมีขั้นตอนวิธีหรืออัลกอริทึมที่รู้จักกันอย่างแพร่หลาย เช่น อัลกอริทึม k-means อัลกอริทึม agglomerative and divisive hierarchical clustering [2] เป็นต้น บทความนี้จะนำอัลกอริทึม covariance-based clustering ที่ผู้เขียนเคยเสนอใน [1] มาใช้ในการวิเคราะห์จัดกลุ่มเอกสารข้อความ โดยจะใช้ข้อมูลตัวอย่างที่กำหนดให้ใน KDD Cup 2002 (Task2) [5] เพื่อแสดงให้เห็นว่าวิธีการจัดกลุ่มดังกล่าวสามารถประยุกต์ใช้งานได้แม้ข้อมูลจะมีขนาดใหญ่

2. การจัดกลุ่มข้อมูลด้วยวิธี Covariance-based clustering

การจัดกลุ่มข้อมูลด้วยวิธี covariance-based clustering [1] อาศัยเทคนิคในการนำเอาแกนเวกเตอร์ หรือสเปกตรัมของเมตริกซ์มาช่วยในการจัดรูปลักษณะของเมตริกซ์ใหม่ โดยสลับแถวและสลับกึ่งของเมตริกซ์ข้อมูล ให้สมาชิกที่มีค่าใกล้เคียงกันปรากฏในบริเวณใกล้เคียงกันมากที่สุด การเปลี่ยนรูปลักษณะของเมตริกซ์ดังกล่าว ทำให้

สามารถเห็นถึงการจัดกลุ่มของสมาชิก ซึ่งหมายถึงการจัดกลุ่มของข้อมูลแต่ละแถวและการจัดกลุ่มของแอททริบิวต์ในแต่ละสดมภ์นั่นเอง

ลำดับการจัดเรียงดังกล่าว ได้จากการวิเคราะห์สเปกตรัลของเมตริกซ์โควาเรนซ์ C (covariance matrix) ของเมตริกซ์ข้อมูล U ที่พิจารณา โดยลำดับการจัดเรียงสดมภ์ π_c ได้จากลำดับการเรียงค่าสมาชิกของเอกเนเวคเตอร์ v_n ที่สอดคล้องกับค่าเอกเนเวคเตอร์ที่มากที่สุดของเมตริกซ์โควาเรนซ์ และ ลำดับการจัดเรียงแถว π_r ได้จากลำดับการเรียงค่าสมาชิกของเวคเตอร์ผลคูณ UV_n ระหว่างเมตริกซ์ข้อมูล U ที่พิจารณากับเอกเนเวคเตอร์ดังกล่าว ในส่วนย่อยต่อไปนี้จะอธิบายที่มาของขั้นตอนวิธีข้างต้น

2.1 การจัดกลุ่มสมาชิกของเมตริกซ์ข้อมูลด้วยวิธีสเปกตรัล

กำหนดให้ U เป็นเมตริกซ์ของจำนวนจริงขนาด $m \times n$ การจัดกลุ่มสมาชิกของ U หมายถึงการจัดลำดับแถวและสดมภ์ของ U ใหม่ ให้สมาชิกของ U อยู่ใกล้เคียงกันมากที่สุด นั่นคือเราสามารถนิยามปัญหาการจัดกลุ่มสมาชิกของเมตริกซ์ข้อมูล U ได้ดังนี้

$$\max_{\pi_c} \frac{1}{m} \sum_{j=1}^{n-1} \langle U_j^{\pi_c}, U_{j+1}^{\pi_c} \rangle + \max_{\pi_r} \frac{1}{n} \sum_{i=1}^{m-1} \langle U_i^{\pi_r}, U_{i+1}^{\pi_r} \rangle$$

โดยที่ พจน์ซ้ายมือ หมายถึง การหาลำดับสดมภ์ที่ทำให้สมาชิกมีการจัดกลุ่ม

พจน์ขวามือ หมายถึง การหาลำดับแถวที่ทำให้สมาชิกมีการจัดกลุ่ม

π_c, π_r หมายถึง เวกเตอร์การจัดเรียงลำดับของสดมภ์และแถวตามลำดับ

$U_j^{\pi_c}$ หมายถึง สดมภ์ที่ j ของ U หลังการจัดเรียงสดมภ์ด้วย π_c

$U_i^{\pi_r}$ หมายถึง แถวที่ i ของ U หลังการจัดเรียงแถวด้วย π_r

$\langle u, v \rangle$ หมายถึง ความคล้ายคลึงกันของเวคเตอร์จำนวนจริง $u, v \in R^p$ ซึ่งนิยามโดย

$$\langle u, v \rangle = \sum_{k=1}^p u_k v_k$$

จะเห็นว่าปัญหาการจัดกลุ่มสมาชิกดังกล่าวสามารถแยกพิจารณาอย่างเป็นอิสระต่อกัน และแต่ละปัญหาย่อยเป็นปัญหา quadratic

combinatorial optimization ซึ่งจัดได้ว่ามีความซับซ้อนของปัญหาอยู่ในกลุ่ม NP-Complete [3] ดังนั้น การหาเวคเตอร์ลำดับที่ดีที่สุดในเวลา polynomial time จึงทำได้ยาก บทความนี้ขอเสนอวิธีสเปกตรัลที่นำเอกเนเวคเตอร์ของเมตริกซ์มาใช้เพื่อหาเวคเตอร์ลำดับ π_c, π_r ที่ให้การจัดกลุ่มของสมาชิกได้ดีภายในเวลาอันสมควร โดยจะแยกพิจารณาการจัดลำดับสดมภ์ในส่วนที่ 2.2 และพิจารณาการจัดลำดับแถวในส่วนที่ 2.3

2.2 การจัดลำดับสดมภ์เพื่อการจัดกลุ่มสมาชิกของเมตริกซ์ข้อมูล

เมื่อกำหนดเมตริกซ์ของจำนวนจริงแบบไม่สมมาตร U ขนาด $m \times n$ มาให้ เราสามารถนิยามปัญหาการจัดลำดับสดมภ์เพื่อการจัดกลุ่มสมาชิกของเมตริกซ์ข้อมูลได้ ดังนี้

$$\max_{\pi_c} \frac{1}{m} \sum_{j=1}^{n-1} \langle U_j^{\pi_c}, U_{j+1}^{\pi_c} \rangle$$

เราสามารถหาความคล้ายคลึงกันของสมาชิกในสดมภ์ $U_j, U_{j'}$ ได้ด้วยค่าผลคูณเวคเตอร์

$$\langle U_j, U_{j'} \rangle = \sum_{i=1}^m U_{i,j} U_{i,j'}$$

ดังนั้นเราสามารถหาความคล้ายคลึงกันเฉลี่ยระหว่างสดมภ์ใดๆด้วยเมตริกซ์ C โดยที่

$$C_{j,j'} = \frac{\langle U_j, U_{j'} \rangle}{m} \quad \forall j, j' \in \{1, \dots, n\}$$

นั่นคือ

$$C = \frac{1}{m} U^T U$$

จะสังเกตได้ว่า C เป็นเมตริกซ์ของจำนวนจริงแบบสมมาตร ค่า $C_{j,j'}$ จะบ่งบอกถึงความคล้ายคลึงกันระหว่างสดมภ์ j และ j' ในเมตริกซ์ U ดังนั้นถ้าต้องการให้สดมภ์ในเมตริกซ์ U ที่มีลักษณะคล้ายกันมาอยู่ใกล้กันจะต้องจัดลำดับเมตริกซ์ C ให้สมาชิกบริเวณเส้นทแยงมุมมีค่าให้มากที่สุด นั่นคือ การจัดกลุ่มสมาชิกของ U โดยการลำดับสดมภ์ใหม่สามารถทำได้โดยการหาลำดับของการจัดเรียงแถวและสดมภ์ของ C ใหม่ ให้สมาชิกที่มีค่ามากอยู่บริเวณเส้นทแยงมุมของเมตริกซ์ C ให้ได้มากที่สุด ปัญหาดังกล่าวมีนิยามที่เทียบเท่าได้ดังนี้

$$\max_{\pi_c} (\pi_c^T C \pi_c) = \max_{\pi_c} \frac{1}{m} \sum_{j=1}^{n-1} \langle U_j^{\pi_c}, U_{j+1}^{\pi_c} \rangle$$

หากพิจารณาค่าเอกเอนและเอกเอนเวกเตอร์ของ C แล้วจะเกิดสมการดังนี้

$$Cv_i = \lambda_i v_i$$

โดยที่ λ_i คือ ค่าเอกเอนที่ i ของเมตริกซ์ C
 v_i คือ เอกเอนเวกเตอร์ที่สอดคล้องกับค่า λ_i โดยที่ $\|v_i\| = 1$

หากคูณสมการข้างต้นด้วย v_i^T ทั้งสองข้างจะได้

$$v_i^T C v_i = v_i^T \lambda_i v_i = \lambda_i$$

หากกำหนดให้ค่าเอกเอนของเมตริกซ์ C เรียงลำดับดังนี้

$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

เอกเอนเวกเตอร์ v_n ที่สอดคล้องกับค่าเอกเอน λ_n ที่มากที่สุด จะเป็นคำตอบของปัญหา $\max_{x \in \mathbb{R}^n} x^T C x$ โดย x เป็นเวกเตอร์ที่มีค่าสมาชิกเป็นจำนวนจริง

อย่างไรก็ตาม ในปัญหาของการจัดเรียงลำดับสดมภ์ของเมตริกซ์ C ที่พิจารณา $\left[\max_{\pi_c} (\pi_c^T C \pi_c) \right]$ เราต้องการ เวกเตอร์ของการจัดลำดับที่มีค่าสมาชิกเป็นจำนวนเต็มบวกที่อยู่ในช่วง 1 ถึง n แบบไม่ซ้ำ ใน [1] มีการเสนอให้ใช้วิธีการ approximation เอกเอนเวกเตอร์ v_n ให้เป็น π_c โดยใช้ลำดับของการจัดเรียงค่าสมาชิกของ v_n เป็นคำตอบ นั่นคือ $\pi_c = \text{order}(v_n)$

วิธีสเปกตรัลในการจัดลำดับสดมภ์เพื่อการจัดกลุ่มสมาชิกของเมตริกซ์ข้อมูล จึงมีขั้นตอนวิธีดังต่อไปนี้

Algorithm CovarianceBasedColumnClustering(U, π_c)

Input: $U \in \mathbb{R}^{m \times n}$ เมตริกซ์ของจำนวนจริงแบบไม่สมมาตร

Output: π_c เวกเตอร์การจัดเรียงลำดับสดมภ์ของเมตริกซ์ U เพื่อให้สมาชิกของเมตริกซ์มีการจัดกลุ่ม

(1) คำนวณหา โควาเรนซ์เมตริกซ์ของเมตริกซ์ข้อมูล

$$C = \frac{1}{m} U^T U$$

(2) คำนวณหาเอกเอนเวกเตอร์ที่สอดคล้องกับ ค่าเอกเอนเวกเตอร์ที่มากที่สุด ของเมตริกซ์ โควาเรนซ์ $v_n = \text{eigenvector}(C, n)$

(3) คำนวณหาลำดับการจัดเรียงสดมภ์ของเมตริกซ์ข้อมูล

$$\pi_c = \text{order}(v_n)$$

2.3 การจัดลำดับแถวเพื่อการจัดกลุ่มสมาชิกของเมตริกซ์ข้อมูล

เมื่อกำหนด U เป็นเมตริกซ์ของจำนวนจริงแบบไม่สมมาตรขนาด $m \times n$ ปัญหาการจัดลำดับแถวเพื่อการจัดกลุ่มสมาชิกของเมตริกซ์ข้อมูล U สามารถนิยามได้ดังนี้

$$\max_{\pi_r} \frac{1}{n} \sum_{i=1}^{m-1} \langle U_{r_i}^{\pi_r}, U_{i+1}^{\pi_r} \rangle$$

โดยที่ π_r เป็นเวกเตอร์ของการจัดลำดับแถว

เราสามารถพิจารณาปัญหาดังกล่าวในลักษณะเดียวกับการพิจารณาใน ส่วน 2.2 นั่นคือ หากความคล้ายคลึงกันเฉลี่ยระหว่างแถว $i, i' \in \{1, \dots, m\}$ ใน U กำหนดได้โดยเมตริกซ์ R โดยที่ $R_{i,i'} = \frac{\langle U_i, U_{i'} \rangle}{n}$ นั่นคือ $R = \frac{1}{n} U U^T$ แล้วเราสามารถนิยามปัญหาการจัดลำดับแถวที่เทียบกันได้ดังนี้

$$\max_{\pi_r} (\pi_r^T R \pi_r) = \max_{\pi_r} \frac{1}{n} \sum_{i=1}^{m-1} \langle U_i^{\pi_r}, U_{i+1}^{\pi_r} \rangle$$

เราสามารถแก้ปัญหาดังกล่าวได้ในลักษณะคล้ายกันกับกรณีการจัดลำดับสดมภ์ นั่นคือการหาเอกเอนเวกเตอร์ที่สอดคล้องกับค่าเอกเอนที่มากที่สุดของ R และหาลำดับการจัดเรียงของสมาชิกของเอกเอนเวกเตอร์ดังกล่าวเพื่อนำไปจัดลำดับแถวของเมตริกซ์ U

อย่างไรก็ตาม โดยส่วนใหญ่ เมตริกซ์ U มักมีจำนวนแถวสูงกว่าจำนวนสดมภ์มาก การคำนวณหาเอกเอนเวกเตอร์ของเมตริกซ์ R จึงไม่สามารถทำได้ง่ายนัก ใน [1] มีการเสนอวิธีการนำเอกเอนเวกเตอร์ของเมตริกซ์ C ซึ่งคำนวณได้จากส่วนที่แล้วมาใช้เพื่ออนุมานหาการจัดลำดับแถวโดยไม่ต้องคำนวณหาเอกเอนเวกเตอร์ของเมตริกซ์ R ใหม่

จะสังเกตได้ว่า $R = \frac{1}{n} U U^T$ และ $C = \frac{1}{m} U^T U$ ดังนั้นเราจะได้ความสัมพันธ์ดังต่อไปนี้

$$U^T R U = \frac{1}{n} U^T U U^T U = \frac{m^2}{n} C C$$

ซึ่งทำให้ปัญหาต่อไปนี้เป็นเทียบเท่ากัน

$$\max_{x \in \mathbb{R}^m} (x^T U^T R U x) = \frac{m^2}{n} \max_{x \in \mathbb{R}^m} (x^T C^2 x)$$

สำหรับทุกๆ เวกเตอร์ v_i ที่สอดคล้องกับค่าเอเจน λ_i ของเมตริกซ์ C เราจะได้

$$C v_i = \lambda_i v_i$$

หากคูณเมตริกซ์ C ทั้งสองข้างของสมการข้างต้น จะทำให้ได้

$$C^2 v_i = \lambda_i C v_i = \lambda_i^2 v_i$$

จากสมการข้างต้น ทำให้เราทราบว่า เมื่อ λ_n เป็นค่าเอเจนที่มากที่สุดสำหรับเมตริกซ์ C แล้ว λ_n^2 จะเป็นค่าเอเจนที่มากที่สุดสำหรับเมตริกซ์ C^2 ด้วย ดังนั้น v_n จึงเป็นเวกเตอร์ที่สอดคล้องกับค่าเอเจนที่มากที่สุดสำหรับเมตริกซ์ C^2

นั่นหมายความว่า v_n เป็นคำตอบของปัญหา $\max_{x \in \mathbb{R}^n} x^T C^2 x$ และปัญหา $\max_{x \in \mathbb{R}^n} x^T U^T R U x$ นั่นคือ $U v_n$ เป็นคำตอบของปัญหา $\max_{x \in \mathbb{R}^m} x^T R x$

จากความจริงดังกล่าว เราสามารถอนุมานได้ว่า เวกเตอร์ของการจัดลำดับแถว π_r สามารถหาได้จากลำดับการจัดเรียงค่าสมาชิกของเวกเตอร์ $U v_n$ โดยที่ v_n เป็นเวกเตอร์ที่สอดคล้องกับค่าเอเจนที่มากที่สุดของเมตริกซ์ C เมื่อรวมการพิจารณาทั้งส่วนที่ 2.2 และ 2.3 ขั้นตอนวิธีสเปกตรัลของการหาการจัดลำดับสมาชิกและแถวของเมตริกซ์ U เพื่อการจัดกลุ่มสมาชิก จึงมีดังนี้

Algorithm CovarianceBasedClustering (U, π_c, π_r)

Input: $U \in \mathbb{R}^{m \times n}$ เมตริกซ์ข้อมูลที่เป็นจำนวนจริงแบบไม่สมมาตร

Output: π_c, π_r เวกเตอร์การจัดเรียงลำดับสมาชิกและแถวของเมตริกซ์ U เพื่อให้สมาชิกเมตริกซ์มีการจัดกลุ่ม

- (1) คำนวณหาโควาเรนซ์เมตริกซ์ของเมตริกซ์ข้อมูล $C = \frac{1}{m} U^T U$
- (2) คำนวณหาเวกเตอร์ที่สอดคล้องกับค่าเอเจนเวกเตอร์ที่มากที่สุด ของเมตริกซ์โควาเรนซ์ $v_n = \text{eigenvector}(C, n)$
- (3) คำนวณหาลำดับการจัดเรียงสมาชิกของเมตริกซ์ข้อมูล $\pi_c = \text{order}(v_n)$

- (4) คำนวณหาลำดับการจัดเรียงแถวของเมตริกซ์ข้อมูล $\pi_r = \text{order}(U v_n)$

การคำนวณหาเอเจนเวกเตอร์ในขั้นตอนที่ (2) เป็นการคำนวณหาเอเจนเวกเตอร์เฉพาะของเมตริกซ์จำนวนจริงแบบสมมาตร ซึ่งสามารถทำได้โดยวิธีทั่วไป เช่น Lanczos's method, Re-orthogonalization, etc [4].

จะสังเกตได้ว่าความซับซ้อนของขั้นตอนวิธีที่เสนอ ขึ้นอยู่กับความซับซ้อนของการคำนวณในแต่ละขั้นตอน โดยสามารถประมาณอย่างคร่าว ๆ ได้ตามลำดับดังนี้ $O(mn^2)$, $O(n^3)$, $O(n \cdot \log n)$, และ $O(mn + n \cdot \log n)$ จะเห็นได้ว่า หาก $m \gg n$ ความซับซ้อนทั้งหมดของวิธีที่เสนอจะเท่ากับ $O(mn^2)$ และหาก $m \ll n$ ความซับซ้อนจะเท่ากับ $O(n^3)$

3. ข้อมูลที่กำหนดให้ใน KDD Cup 2002

ในส่วนนี้เราจะนำเทคนิคการจัดกลุ่มข้อมูลด้วยวิธี covariance-based clustering มาประยุกต์ใช้ในการศึกษาลักษณะของข้อมูลชุดเอกสารต่างๆ ที่กำหนดให้ใน KDD Cup 2002 (Task2) [5] ซึ่งมีส่วนประกอบต่างๆ ออกเป็นดังนี้

3.1 Gene-Abstracts file

Gene-Abstracts เป็นแฟ้มข้อมูลที่ระบุให้ทราบว่ายีน (gene) แต่ละยีนถูกกล่าวถึงในบทคัดย่อ (abstracts) ใดบ้าง โดยในบทคัดย่อจะมีคำสำคัญ (keyword) อธิบายคุณลักษณะต่างๆ ที่เกี่ยวกับยีนนั้นภายในแฟ้มข้อมูล Gene-Abstracts จะมีลักษณะดังรูปที่ 1

YML034W	10734188
YML034W	10894548
YHR051W	207698
:	:

รูปที่ 1 แสดงตัวอย่างข้อมูลภายในแฟ้ม Gene-Abstracts

ในแต่ละบรรทัดภายในแฟ้มข้อมูลแสดงความสัมพันธ์เชิงคู่ลำดับระหว่างรหัสของยีน(สมาชิกแรก) กับชื่อแฟ้มข้อมูลของบทคัดย่อ(สมาชิกที่สอง) ยีนที่พิจารณาเป็นยีนของเซตยีสต์โดยระบุไว้ในที่นี้ทั้งสิ้น 3329 ยีน

3.2 Abstracts Files

บทความทั้งหมด 15,234 บทความ ข้อมูล อยู่ในไฟล์เดียวกัน แต่ละบทความถูกกำหนดชื่อด้วยรหัสตัวเลขซึ่งสอดคล้องกับที่อ้างอิงในแฟ้มข้อมูล Gene-Abstracts บทความเหล่านี้ สามารถดาวน์โหลดได้จากฐานข้อมูล MEDLINE database (www.ncbi.nlm.nih.gov/entrez/query.fcgi) ภายในแฟ้มข้อมูลของบทความจะมีลักษณะเป็นดังรูปที่ 2

Some strains of Saccharomyces cerevisiae exhibit a specific transport system for ureidosuccinic acid, which is regulated by nitrogen metabolism. ...

รูปที่ 2 แสดงตัวอย่างข้อมูลภายในแฟ้มบทความเลขที่ 13831

โดยอาจมีการปรับแฟ้มข้อมูล Function-Hierarchy ใหม่ เพื่อให้เหมาะสมกับการทำหน้าที่เป็นพจนานุกรมคำศัพท์ (dictionary) ด้วยการคัดคำบางคำที่ไม่สนใจหรือไม่สำคัญออกจากแฟ้มเอกสารนี้ ภายในแฟ้มข้อมูล Function-Hierarchy ใหม่จะมีลักษณะดังตัวอย่างในรูปที่ 3 ซึ่งจะได้ฟังก์ชันของยีนทั้งสิ้น 342 ฟังก์ชัน

METABOLISM
amino acid metabolism
amino acid biosynthesis
biosynthesis aspartate
biosynthesis lysine
biosynthesis cysteine-aromatic

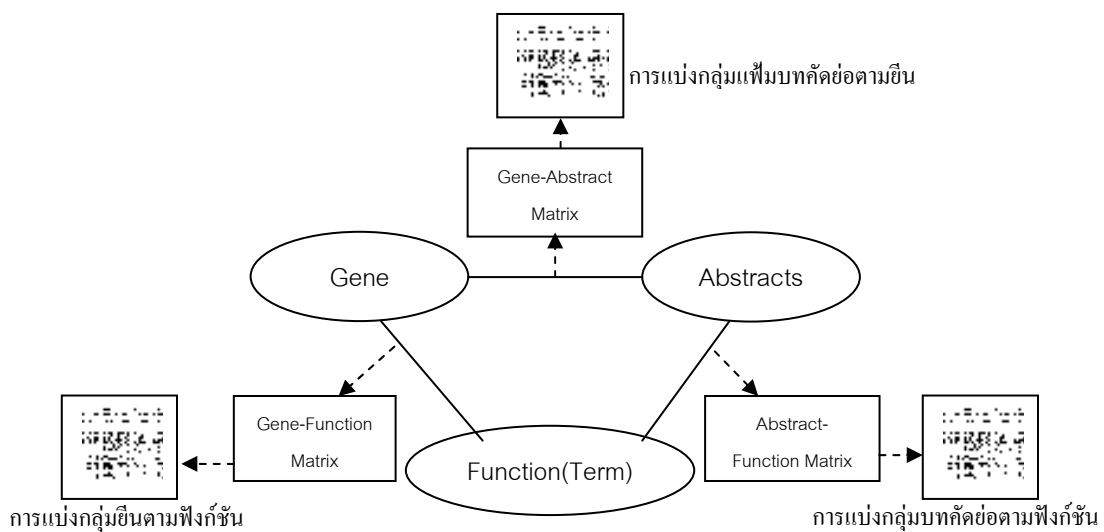
รูปที่ 3 แสดงตัวอย่างข้อมูลภายในแฟ้ม Function-Hierarchy

3.3 Function-Hierarchy file

ในแฟ้มข้อมูลนี้เป็นการอธิบายข้อมูลในลักษณะลำดับชั้น (hierarchy) ของฟังก์ชันต่างๆ (function) ของยีน ในการจัดกลุ่มเอกสารบทความ (abstract clustering) นั้น เราสามารถนำ Function-Hierarchy มาใช้สำหรับกำหนดคำศัพท์ทางเทคนิค (technical term)

4. การวิเคราะห์จัดกลุ่มข้อมูล KDD Cup 2002

การจัดกลุ่มเอกสารงาน KDD Cup Task2 สามารถนำขั้นตอนต่างๆ มาเขียนเป็นแผนผังได้ดังรูปที่ 4



รูปที่ 4 แสดงแผนผังขั้นตอนการวิเคราะห์การจัดกลุ่มเอกสาร KDD Cup 2002

4.1 การจัดกลุ่มของยีนและเพิ่มบทคัดย่อ

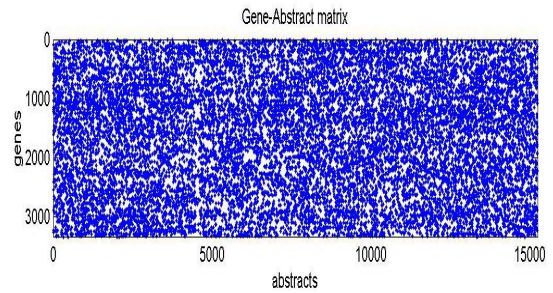
ขั้นตอนการจัดกลุ่มข้อมูลเริ่มจากการนำเพิ่มข้อมูล Gene-Abstracts มาจัดรูปแบบใหม่ให้เหมาะสมเพื่อนำไปสร้างเป็น เมตริกซ์สำหรับการนำไปจัดกลุ่ม โดยทำการกำหนดบรรทัดนี้ให้กับยีนและเพิ่มบทคัดย่อ ซึ่งจะได้อันดับของความสัมพันธ์ระหว่างยีนกับเพิ่มบทคัดย่อ ทำให้ได้รูปแบบดังตัวอย่างนี้

[3306]	'YPR173C'	'10393249'	[604]
[3306]	'YPR173C'	'1097413'	[2561]
[3306]	'YPR173C'	'11329380'	[3568]
[3306]	'YPR173C'	'11559748'	[4160]
[3306]	'YPR173C'	'11563910'	[4176]
[3306]	'YPR173C'	'9155008'	[12497]
[3306]	'YPR173C'	'9431454'	[13546]

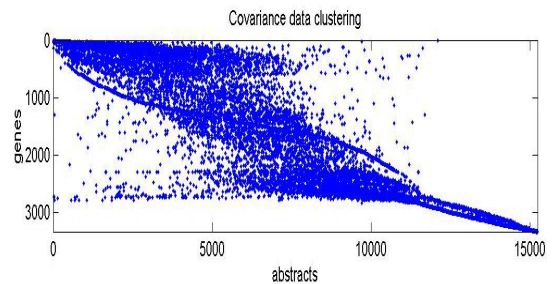
ในสคริปต์แรกจะเป็นบรรทัดของยีน (gene index) และสคริปต์สุดท้ายเป็นบรรทัดของเพิ่มบทคัดย่อ (abstract index) เมื่อมีการกำหนดบรรทัดนี้เรียบร้อยแล้วจึงสร้างเป็นเมตริกซ์ ซึ่งในที่นี้เราใช้เมตริกซ์แบบ sparse matrix ที่สามารถสร้างเมตริกซ์ขนาดใหญ่และรองรับข้อมูลมากๆ ได้เนื่องจากจำนวนของชื่อเพิ่มข้อมูลของบทคัดย่อมีจำนวนมากถึง 15,234 เพิ่มข้อมูล

หากใช้การจัดการเมตริกซ์ในรูปแบบปกติ จะทำให้เกิดปัญหาในเรื่องของหน่วยความจำและเวลาในการทำงานของโปรแกรม เนื่องจากเมตริกซ์แบบเต็มรูปของความสัมพันธ์ระหว่างยีนกับข้อมูลบทคัดย่อจะมีขนาดเป็นจำนวนของยีนคูณกับจำนวนของเพิ่มบทคัดย่อ ข้อดีของการใช้ sparse matrix คือ สามารถจัดการข้อมูลในลักษณะคู่อันดับ ทำให้สามารถนำบรรทัดที่กำหนดมาระบุลงในเมตริกซ์รูปแบบนี้ได้ทันทีและสะดวกกับการจัดการ

จากความสัมพันธ์ของยีนกับเพิ่มบทคัดย่อข้างต้น เราสามารถวิเคราะห์หาการจัดกลุ่มของยีนที่ถูกกล่าวถึงในเพิ่มบทคัดย่อด้วยการจัดกลุ่มข้อมูลโดยวิธีการ covariance-based clustering ผลลัพธ์แสดงดังรูปที่ 5 และ 6 โดยในรูปที่ 5 และ 6 แกนนอนจะหมายถึงเพิ่มบทคัดย่อ แกนตั้งจะหมายถึงยีน จากรูปเราสามารถวิเคราะห์ได้อย่างคร่าวๆ ว่ายีนแบ่งออกเป็นสองกลุ่มใหญ่ๆ โดย กลุ่มหนึ่งมีจำนวนบทคัดย่อที่เกี่ยวข้องจำนวนน้อย และอีกกลุ่มหนึ่งมีจำนวนบทคัดย่อที่เกี่ยวข้องจำนวนมาก เราสามารถจำแนกรายชื่อของยีนและเพิ่มบทคัดย่อของทั้งสองกลุ่มดังกล่าวเพื่อประโยชน์ในการสืบค้นต่อไป



รูปที่ 5 แสดงแผนภาพก่อนการจัดกลุ่มยีนกับบทคัดย่อ

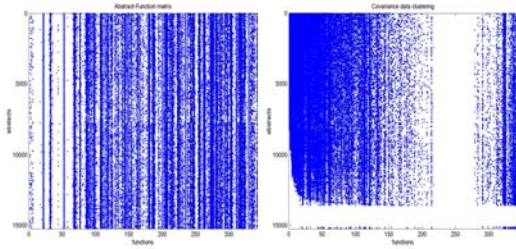


รูปที่ 6 แสดงแผนภาพหลังจากการจัดกลุ่มยีนกับบทคัดย่อ

4.2 การจัดกลุ่มของฟังก์ชันของยีนและเพิ่มบทคัดย่อ

ในส่วนของการทำงานการจัดกลุ่มของบทคัดย่อตามความคล้ายของจำนวนฟังก์ชันที่ปรากฏอยู่ภายในบทคัดย่อสามารถทำได้โดยการสร้างเมตริกซ์ความสัมพันธ์ในรูปแบบ Abstracts-Function คำศัพท์ที่สำคัญในแต่ละเพิ่มบทคัดย่อหมายถึงฟังก์ชันที่ระบุไว้ในเพิ่มข้อมูล Function-Hierarchy ที่ได้กล่าวข้างต้น เราจะนำคำศัพท์ชุดนี้มาทำการคัดเลือกเฉพาะคำที่สำคัญออกจากบทคัดย่อ โดยวนเปิดเพิ่มบทคัดย่อต่างๆ แล้ว ตรวจสอบฟังก์ชัน จากนั้นจะได้เป็นเมตริกซ์ ที่มีขนาดเป็นจำนวนของเพิ่มบทคัดย่อคูณกับจำนวนของฟังก์ชัน โดยสมาชิกของเมตริกซ์นี้ จะระบุจำนวนครั้งที่ปรากฏของฟังก์ชันในแต่ละเพิ่มบทคัดย่อ

เราสามารถวิเคราะห์หาความสัมพันธ์ของฟังก์ชันในแต่ละเพิ่มข้อมูล โดยใช้วิธีการจัดกลุ่มข้อมูลแบบ covariance-based clustering อีกครั้ง ซึ่งผลลัพธ์เป็นดังรูปที่ 7 และทำให้สามารถสรุปอย่างคร่าวๆ ได้ว่าบทคัดย่อส่วนใหญ่กล่าวถึงฟังก์ชันต่างๆ ของยีน และจะมีบางฟังก์ชันที่ไม่ได้ถูกกล่าวถึงเลย

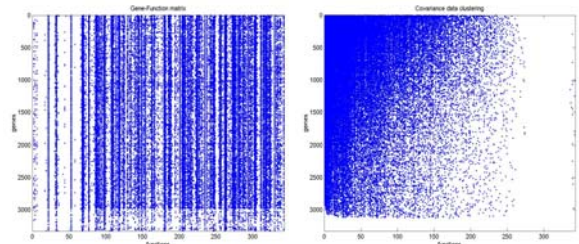


รูปที่ 7 แสดงแผนภาพก่อนการทำการจัดกลุ่มและหลังจากการจัดกลุ่มคำศัพท์เฉพาะในแต่ละแฟ้มบทคัดย่อ โดยแกนตั้งหมายถึงแฟ้มบทคัดย่อ แกนนอนจะหมายถึงฟังก์ชันของยีน

4.3 การจัดกลุ่มของยีนและฟังก์ชันของยีน

ในการศึกษาทางชีววิทยาระดับโมเลกุล เราต้องการทราบว่าแต่ละยีนของเซลล์มีชีวิตมีฟังก์ชันอย่างไรบ้าง ในส่วนนี้เราจะนำผลจากการวิเคราะห์ข้อมูล ในสองส่วนข้างต้นมาวิเคราะห์หาการจัดกลุ่มระหว่างยีนและฟังก์ชันของยีน จากเมตริกซ์ยีนกับบทคัดย่อ (Gene-Abstracts) และ เมตริกซ์คำศัพท์เฉพาะที่สำคัญในแต่ละแฟ้มบทคัดย่อ (Abstracts-Function) เมื่อเราต้องการหาความสัมพันธ์ระหว่าง ยีน กับคำศัพท์ในแต่ละแฟ้มบทคัดย่อ (Function) เราสามารถนำเมตริกซ์ข้างต้นทั้งสองมาคูณกัน เพื่อให้ได้เมตริกซ์ผลลัพธ์ Gene-Function = Gene-Abstracts * Abstracts-Function ค่าสมาชิกของเมตริกซ์ Gene-Function จะหมายถึง จำนวนบทคัดย่อที่กล่าวถึงยีนและฟังก์ชัน ซึ่งบ่งบอกความสัมพันธ์โดยนัยระหว่างยีนและฟังก์ชันนั้นๆ

เมื่อได้เมตริกซ์ Gene-Function เราสามารถใช้วิธี covariance-based clustering ในการจัดกลุ่ม ซึ่งทำให้ได้ผลสรุปออกมาดังรูปที่ 8 ผลจากการจัดกลุ่มในรูปที่ 8 ทำให้สามารถสรุปอย่างคร่าวได้ว่ายีนที่พิจารณามีกลุ่มฟังก์ชันเฉพาะร่วมกันอย่างเรียงราย และมีกลุ่มยีนที่ยังไม่มีการค้นพบหน้าที่หรือฟังก์ชันของมัน



รูปที่ 8 แสดงแผนภาพก่อนการทำการจัดกลุ่มและหลังจากการจัดกลุ่มยีนกับคำศัพท์ในแต่ละแฟ้มบทคัดย่อ โดยแกนนอนหมายถึงฟังก์ชันของยีน แกนตั้งหมายถึงยีน

ขั้นตอนวิธีต่างๆ ที่อธิบายข้างต้นได้รับการพัฒนาด้วยโปรแกรมแมทแลป และปฏิบัติการบนเครื่องคอมพิวเตอร์ส่วนบุคคลที่มีหน่วยประมวลผลรุ่น Pentium III ความเร็ว 650 MHz. และมีหน่วยความจำ 256 Mb. ตารางที่ 1 แสดงขนาดและเวลาการคำนวณของวิธีการจัดกลุ่มสมาชิกของเมตริกซ์ต่างๆ ด้วยวิธีสเปกตรัล

ตารางที่ 1 ขนาดและเวลาการคำนวณของวิธีการจัดกลุ่มสมาชิกของเมตริกซ์ต่างๆ ด้วยวิธีสเปกตรัล

เมตริกซ์	ยีนและแฟ้มบทคัดย่อ	แฟ้มบทคัดย่อและฟังก์ชันของยีน	ยีนและฟังก์ชันของยีน
ขนาด	3329*15234	15234*342	3329*342
เวลาการคำนวณ (วินาที)	1.28	23.22	5.05

5. สรุป

บทความนี้ได้แสดงการนำวิธีการจัดกลุ่มสมาชิกด้วยวิธี covariance-based clustering มาประยุกต์ใช้ในการวิเคราะห์ข้อมูลเอกสารขนาดใหญ่ วิธีการจัดกลุ่มสมาชิกดังกล่าวอาศัยการวิเคราะห์ปริมาณการจัดกลุ่มสมาชิกด้วยค่าเอเจนและเอเจนเวกเตอร์ของเมตริกซ์โควาเรนซ์ของเมตริกซ์ที่พิจารณา จากคุณสมบัติทางคณิตศาสตร์ต่างๆ เกี่ยวกับค่าเอเจนและเอเจนเวกเตอร์ ทำให้เราก็ค้นพบลำดับการจัดเรียงแถวและสดมภ์ เพื่อการจัดกลุ่มสมาชิกได้อย่างมีประสิทธิภาพและสามารถประมวลผลได้กับข้อมูลขนาดใหญ่

เอกสารอ้างอิง

- [1] จิรยุทธ ไชยจรรูวณิช สุวิษยะ รัตตะระมย์ วิธีสเปกตรัลเพื่อการจัดกลุ่มสมาชิกของเมตริกซ์แบบไม่สมมาตร The Sixth National Computer Science and Engineering Conference (NCSEC 2002), October 29-31 2002, Pattaya, Thailand
- [2] J.Han and M.Kamber. Data Mining: Concepts and Techniques. Morgan Kaufmann Publishers. 2001
- [3] A.L.Peressini, F.E.Sullivan, and J.J.Uhl Jr., The Mathematics of Nonlinear Programming, Springer-Verlag, New York, 1980
- [4] J.H.Wilkinson, The algebraic eigenvalue problem, Clarendon press, Oxford, 1965
- [5] KDD Cup 2002, The Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, <http://www.biostat.wisc.edu/~craven/kddcup/>



จิรยุทธ ไชยจรรูวณิช สำเร็จการศึกษา BSc. (Computer Science) ปี 1995, MSc. (Computer Science) ปี 1996 จาก University of Grenoble I ประเทศฝรั่งเศส และ PhD. (Computer Science) ปี 2000 จาก University of Evry Val d'Essonne ประเทศฝรั่งเศส ปัจจุบันเป็นอาจารย์ประจำภาควิชาวิทยาการคอมพิวเตอร์ คณะวิทยาศาสตร์ มหาวิทยาลัยเชียงใหม่