

การเพิ่มความเร็ว SOM ในการจัดกลุ่มข้อมูลด้วยเทคนิคการสุ่มแบบเบี่ยงเบนตามความหนาแน่นและความเป็นปึกแผ่นของข้อมูล

Increasing speed SOM for Clustering Data with A Density-Biased Sampling and Partial Approximate Compactness Estimator

บัณฑิต บุญวัฒน์นะ บุญวัฒน์ อัครชู

ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

บทคัดย่อ

Self-Organizing Map(SOM) เป็นนิเวรอนเน็ตเวิร์กแบบไม่มีผู้สอน จุดเด่นคือเหมาะสำหรับนำมาวิเคราะห์ข้อมูล โดยข้อมูลที่ได้จะอยู่ในรูปแบบของกริดที่แสดงเป็นแผนภาพช่วยในการวิเคราะห์ข้อมูลได้เป็นอย่างดี จุดอ่อนของแผนภาพ SOM ก็คือจะใช้ระยะเวลาในการเรียนรู้นานเมื่อแผนภาพมีขนาดใหญ่และข้อมูลมีปริมาณมาก ซึ่งการลดขนาดของข้อมูลที่ใช้เรียนรู้ของ SOM เป็นอีกแนวทางหนึ่งที่จะช่วยแก้ปัญหานี้ได้ งานวิจัยนี้จึงนำวิธีการลดขนาดข้อมูลด้วยเทคนิคการสุ่มแบบเบี่ยงเบนตามความหนาแน่นและความเป็นปึกแผ่นของข้อมูลซึ่งสามารถสร้างชุดข้อมูลสุ่มแบบต่อเนื่องได้โดยการอ่านข้อมูลเพียงรอบเดียวแล้วนำข้อมูลที่ได้จากการสุ่มไปทำการเรียนรู้ในแผนภาพของ SOM เพื่อให้แผนภาพมีการจัดเรียงตัวได้เร็วขึ้นแล้วจึงนำข้อมูลทั้งหมดมาเรียนรู้อีกครั้ง พบว่าวิธีที่นำเสนอสามารถลดเวลาในการจัดข้อมูลได้มากกว่า 50 เปอร์เซ็นต์เมื่อเปรียบเทียบกับ SOM แบบเดิมและประสิทธิภาพของการจัดกลุ่มเอกสารยังคงเดิม

คำสำคัญ : เซลล์ออร์แกนไนซิงแมพ, การสุ่มแบบเบี่ยงเบนตามความหนาแน่น, ระยะเวลาแบบยูคลิด

Abstract

Self-Organizing Map(SOM) is an unsupervised neural network providing cluster analysis of high dimensional input data. Outputs from SOM are represented in map that help us to explore data. The weak point of conventional SOM is that it takes long time to train system in case of large map and large data set. Data reduction is an important step to increase the efficiency of SOM. This research presents A Density-Biased Sampling and Partial Approximate Compactness Estimator (DBSPACE) to reduce data. This method requires only one pass of dataset. The output then is used for training SOM to facilitate fast arrangement. Lastly, the learning process is repeated with the whole data. The result from experiment shows that this method can reduce time up to 50 percent, compared with the conventional SOM, yet still get the same result to find cluster.

Key words: Self-Organizing Map, A Density-Biased Sampling, Euclidean distance

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1. บทนำ

SOM ถูกนำมาใช้อย่างกว้างขวางในหลายๆด้านเช่น การประยุกต์ใช้ SOM ในการจดจำเสียง[1] การประยุกต์ใช้ SOM ในการจัดกลุ่มข่าวบนอินเทอร์เน็ต[2] เป็นต้น นอกจากนี้แล้วยังมีงานวิจัยอื่นๆเกี่ยวกับ SOM ได้รับการตีพิมพ์มากกว่า 3000 บทความในช่วงปี ค.ศ.1981-1997[3] อย่างไรก็ตามปัญหาสำคัญในการประยุกต์ใช้งาน SOM คือ เวลาในการเรียนรู้ ซึ่งจะใช้เวลาในการเรียนรู้ค่อนข้างนาน ถ้าข้อมูลมีมิติสูง เช่น เวกเตอร์ของเอกสารมีขนาดมากจะยิ่งใช้เวลาในการเรียนรู้นาน จึงมีงานวิจัยที่พัฒนาลดเวลาในการเรียนรู้ เช่น การนำเอาอัลกอริทึม K-mean ช่วยในการจัดกลุ่มข้อมูลก่อนนำไปเรียนรู้อีกในแผนภาพเพื่อให้แผนภาพเรียงตัวได้เร็วขึ้นในจำนวนรอบที่น้อยกว่า[4] การลดการคำนวณโหนดขณะ โดยแบ่งเป็นส่วนย่อยๆ ที่ใช้ในการค้นหาวิธีนี้สามารถลดเวลาในการคำนวณได้สูงสุดถึง 14 เปอร์เซ็นต์[5] งานวิจัยนี้้นำวิธีการลดขนาดข้อมูลด้วยเทคนิคการสุ่มแบบเบี่ยงเบนตามความหนาแน่นและความเป็นปึกแผ่นของข้อมูลแล้วนำข้อมูลที่ได้จากการสุ่มไปทำการเรียนรู้ในแผนภาพของ SOM เพื่อแก้ปัญหาการใช้เวลาในการเรียนรู้นานโดยประสิทธิภาพการจัดกลุ่มยังคงเดิม

2. ทฤษฎีที่เกี่ยวข้อง

2.1 เทคนิคการสุ่มแบบเบี่ยงเบนตามความหนาแน่นและความเป็นปึกแผ่นของข้อมูล (DBSPACE: A Density-Biased Sampling and Partial Approximate Compactness Estimator)

DBSPACE [6]เป็นเทคนิคที่ใช้สำหรับงานการจัดกลุ่มข้อมูลขนาดใหญ่ ลักษณะที่สำคัญของอัลกอริทึมคือสามารถสร้างชุดข้อมูลสุ่มได้ด้วยการอ่านข้อมูลเพียงแค่รอบเดียวเท่านั้น, รองรับข้อมูลที่มีขนาดใหญ่มากและไม่สามารถทราบจำนวนข้อมูลทั้งหมดมาก่อนล่วงหน้าการสุ่มข้อมูล, สร้างชุดข้อมูลสุ่มที่ดีที่สุดสำหรับข้อมูลที่มีการกระจายแบบไม่ปกติ (Zipf), สามารถตรวจสอบและกรองข้อมูลรบกวนไม่ให้ถูกเลือกเข้ามาในชุดข้อมูลสุ่ม โดยใช้ค่าประมาณความเป็นปึกแผ่นของข้อมูลประกอบการสุ่มข้อมูล การคำนวณค่าต่างๆที่จำเป็นสำหรับอัลกอริทึมมีเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่ออนุญาตให้นำไปใช้

ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ดังนั้น ให้ $\{\bar{x}_i\}$ เป็นเซตข้อมูลขนาด d มิติจำนวน N อ็อบเจกต์ของคลัสเตอร์ใดๆ บน data space โดยที่ $i=1,2,\dots,N$

Centroid ($\bar{\mu}$) หรือจุดศูนย์กลางมวลของคลัสเตอร์เป็นจุดที่ใช้เป็นตัวแทนของทุกอ็อบเจกต์ในคลัสเตอร์ หากได้คังสมการที่ (1)

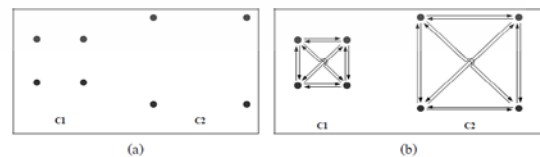
$$\bar{\mu} = \frac{\sum_{i=1}^N \bar{x}_i}{N} \quad (1)$$

Diameter (D) หรือเส้นผ่านศูนย์กลางของคลัสเตอร์เป็นค่าระยะทางเฉลี่ยของทุกอ็อบเจกต์ภายในคลัสเตอร์เดียวกัน ให้ $d(\bar{x}_i, \bar{x}_j)$ เป็นระยะทางตามแบบยูคลิด (Euclidean distance) ระหว่างอ็อบเจกต์ \bar{x}_i กับ \bar{x}_j diameter ของคลัสเตอร์สามารถหาได้จาก

$$D = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N d(\bar{x}_i, \bar{x}_j)^2}{N(N-1)}} \quad (2)$$

$$d(\bar{x}_i, \bar{x}_j) = \sqrt{(\bar{x}_{i1} - \bar{x}_{j1})^2 + (\bar{x}_{i2} - \bar{x}_{j2})^2 + \dots + (\bar{x}_{id} - \bar{x}_{jd})^2} \quad (3)$$

diameter ของคลัสเตอร์แสดงถึงความใกล้ชิดกันระหว่างอ็อบเจกต์ภายในคลัสเตอร์ได้ โดยรูปที่ 1a จะเห็นได้ว่าคลัสเตอร์ C1 มีความหนาแน่นและความเป็นปึกแผ่นของข้อมูลสูงกว่าคลัสเตอร์ C2 และเมื่อลากเส้นเชื่อมระหว่างอ็อบเจกต์ใดๆ ไปยังทุกอ็อบเจกต์ภายในคลัสเตอร์ (ดังรูปที่ 1b) และนำความยาวกำลังสองของเส้นเชื่อมเหล่านั้นมารวมกันแล้วหารด้วยจำนวนเส้นเชื่อมทั้งหมด (12 เท่ากันทั้งสองคลัสเตอร์) และนำมาถอดรากที่สองจะพบว่าค่าที่ได้จากคลัสเตอร์ C1 น้อยกว่าค่าที่ได้จากคลัสเตอร์ C2 แสดงให้เห็นว่าอ็อบเจกต์ของคลัสเตอร์ C1 มีความเป็นปึกแผ่นสูงกว่าอ็อบเจกต์ของคลัสเตอร์ C2



รูปที่ 1 ตัวอย่างของสองคลัสเตอร์ที่มีความหนาแน่นต่างกัน

Compactness (Θ) คือค่าความเป็นปึกแผ่นของข้อมูลแบบประมาณ โดยให้การอ่านข้อมูลเพียงรอบเดียว โดยค่า Compactness เป็นส่วนกลับของ discordancy (Γ)

หรือความไม่ประสานกันภายในกลุ่มข้อมูล โดยค่า Compactness สามารถคำนวณได้ดังสมการที่ (4)

$$\Theta = \frac{1}{\Gamma} \tag{4}$$

สำหรับค่า discordancy สามารถคำนวณได้สองแบบคือ weighted discordancy และ non-weighted discordancy

Weighted discordancy ($\bar{\Gamma}$) คือค่าความไม่ประสานกันภายในกลุ่มข้อมูลแบบให้น้ำหนัก สำหรับการคำนวณค่า $\bar{\Gamma}$ สามารถคำนวณได้ดังสมการที่ (5)

$$\bar{\Gamma} = \frac{\sum_{i=1}^N [(n'_i - 1) \cdot d(\bar{\mu}_{i-1}, \bar{x}_i)]}{\sum_{i=1}^N (n'_i - 1)} \tag{5}$$

โดยให้ $\bar{\mu}_0 = \bar{x}_i$, n'_i เป็นจำนวนอ็อบเจกต์อ้างอิงของตัวแทน และ $n'_i = 2$

Non-weighted discordancy ($\tilde{\Gamma}$) คือค่าความไม่ประสานกันภายในกลุ่มข้อมูลแบบไม่ให้น้ำหนัก โดยค่า $\tilde{\Gamma}$ สามารถคำนวณได้ดังสมการที่ (6)

$$\tilde{\Gamma} = \frac{\sum_{i=1}^N d(\bar{\mu}_{i-1}, x_i)}{(N-1)} \tag{6}$$

ตารางที่ 1 แสดงการคำนวณค่า $\bar{\Gamma}$, $\tilde{\Gamma}$ และ Centroid ($\bar{\mu}$) ของอ็อบเจกต์ \bar{x}_i

i	$\bar{\Gamma}$	$\tilde{\Gamma}$	$\bar{\mu}$
1	$0 \cdot d(\bar{\mu}_0, \bar{x}_1)$	$d(\bar{\mu}_0, \bar{x}_1)$	$\frac{\bar{x}_1}{1}$
2	$\frac{0 \cdot d(\bar{\mu}_0, \bar{x}_1) + 1 \cdot d(\bar{\mu}_1, \bar{x}_2)}{1}$	$\frac{d(\bar{\mu}_0, \bar{x}_1) + d(\bar{\mu}_1, \bar{x}_2)}{1}$	$\frac{\bar{x}_1 + \bar{x}_2}{2}$
3	$\frac{\{0 \cdot d(\bar{\mu}_0, \bar{x}_1) + 1 \cdot d(\bar{\mu}_1, \bar{x}_2) + 2 \cdot d(\bar{\mu}_2, \bar{x}_3)\}}{1+2}$	$\frac{\{d(\bar{\mu}_0, \bar{x}_1) + d(\bar{\mu}_1, \bar{x}_2) + d(\bar{\mu}_2, \bar{x}_3)\}}{2}$	$\frac{\bar{x}_1 + \bar{x}_2 + \bar{x}_3}{2}$
⋮	⋮	⋮	⋮
N	$\bar{\Gamma} = \frac{\sum_{i=1}^N [(n'_i - 1) \cdot d(\bar{\mu}_{i-1}, \bar{x}_i)]}{\sum_{i=1}^N (n'_i - 1)}$	$\tilde{\Gamma} = \frac{\sum_{i=1}^N d(\bar{\mu}_{i-1}, x_i)}{(N-1)}$	$\bar{\mu} = \frac{\sum_{i=1}^N \bar{x}_i}{N}$

สมมุติให้ข้อมูลขนาด N อ็อบเจกต์ x_1, x_2, \dots, x_N ซึ่งถูกแบ่งออกเป็นคลัสเตอร์หรือกลุ่มย่อยๆ g กลุ่ม โดยแต่ละกลุ่มประกอบด้วยอ็อบเจกต์จำนวน n_1, n_2, \dots, n_g (จำนวนของอ็อบเจกต์แสดงถึงความหนาแน่นของข้อมูลภายในกลุ่ม) และความเป็นปึกแผ่นของข้อมูลภายในกลุ่ม (Compactness) มีค่าเท่ากับ $\Theta_1, \Theta_2, \dots, \Theta_g$ เมื่อต้องการสร้าง

ชุดข้อมูลสุ่มขนาด M อ็อบเจกต์ โดยที่ความน่าจะเป็นที่อ็อบเจกต์ x_i ใดๆ ภายในกลุ่มที่ i จะถูกเลือกจะต้อง

(i) ผกผันกับค่าความหนาแน่นของข้อมูลภายในกลุ่ม

$$P(x_i) \propto \frac{1}{n_i} \tag{7}$$

(ii) แปรผันตรงกับค่าความเป็นปึกแผ่นของข้อมูลภายในกลุ่ม

$$P(x_i) \propto \Theta_i \text{ หรือ } P(x_i) \propto \frac{1}{\Gamma_i} \tag{8}$$

ดังนั้นเพื่อให้สอดคล้องกับ (i) และ (ii) จึงให้ความน่าจะเป็นที่อ็อบเจกต์ x_i ใดๆ ภายในกลุ่มที่ i จะถูกเลือกมีค่าเท่ากับ

$$P(x_i) = \frac{\alpha}{n_i^e \cdot \Gamma_i^\delta} \tag{9}$$

โดยให้ α, e , และ δ เป็นค่าคงที่ใดๆ ซึ่งค่า α ที่เหมาะสมสำหรับการสุ่มที่ต้องการชุดข้อมูลสุ่มขนาด M อ็อบเจกต์สามารถพิจารณาได้ดังสมการที่ (10)

$$\alpha = \frac{M}{\sum_{i=1}^g [n_i^{1-e} \cdot \Gamma_i^{-\delta}]} \tag{10}$$

และเมื่อแทนค่า α ในสมการที่ (3)-(4) จะให้ความน่าจะเป็นที่อ็อบเจกต์ x_i ใดๆ ภายในกลุ่มที่ i จะถูกเลือกสำหรับชุดข้อมูลสุ่มขนาด M มีค่าเท่ากับ

$$P(x_i) = \frac{M}{n_i^e \cdot \Gamma_i^\delta \cdot \sum_{i=1}^g [n_i^{1-e} \cdot \Gamma_i^{-\delta}]} \tag{11}$$

จากสมการข้างต้นจะเห็นได้ว่าในบางกรณีส่วนหารสามารถมีค่าเป็นศูนย์ได้ เช่น n_i มีค่าเท่ากับ 1 ซึ่งทำให้ค่า Γ_i มีค่าเป็นศูนย์ ดังนั้นเพื่อเป็นการเลี่ยงปัญหาการหารด้วยค่าศูนย์ (divide by zero) จึงปรับปรุงฟังก์ชัน $d(\bar{x}_i, \bar{x}_i)$ ที่ใช้สำหรับอัลกอริทึม DBSPACE โดยการบวก 1 เพิ่มเข้าไปกับค่าของระยะห่างที่ได้จากการคำนวณปกติ ซึ่งจะมีผลทำให้ค่า Γ ที่ได้มีค่าเพิ่มขึ้นจากเดิมประมาณ 1 มีผลทำให้ Γ^δ เป็นฟังก์ชันเพิ่มเสมอ ในรูปที่ 2 แสดงอัลกอริทึม DBSPACE โดยที่ $n[i]$ แสดงถึงจำนวนข้อมูลภายในกลุ่มที่ i , $lsx[i]$ เป็นผลรวมเชิงเวกเตอร์ของทุกอ็อบเจกต์ภายในกลุ่มที่ i , $lsd[i]$ คือผลรวมของระยะห่างระหว่างอ็อบเจกต์ภายในกลุ่มที่ i

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

 $\alpha_D = 0$ 
FOR each input point  $x$  DO
   $i = \text{hash}( \bar{X} )$ 
  IF  $n[i] < 0$  THEN
 $\alpha_D = \alpha_D - (n[i]^{1-e} \cdot \text{discordancy}(\text{lsd}[i], n[i]))^{-\delta}$ 
 $n[i] = n[i] + 1$ 
  IF  $n[i] = 1$  THEN  $\bar{\mu} = \bar{x}$ 
  IF  $\Gamma = \Gamma$  THEN  $\text{lsd}[i] = \text{lsd}[i] + (n[i]-1) \cdot d(\bar{\mu}, \bar{x})$ 
  ELSE  $\text{lsd}[i] = \text{lsd}[i] + d(\bar{\mu}, \bar{x})$ 
   $\text{ls}\bar{x} = \text{ls}\bar{x} + \bar{x}$ 
   $\bar{\mu} = \text{ls}\bar{x} / n[i]$ 
   $\Gamma = \text{discordancy}(\text{lsd}[i], n[i])$ 
 $\alpha_D = \alpha_D + (n[i]^{1-e} \cdot \Gamma^{-\delta})$ 
WITH prob.  $P = \min \{ M / [\alpha_D \cdot n[i]^e \cdot \Gamma^\delta], 1 \}$  DO
  IF the output buffer is full THEN reduce()
  // start drawing when group contains more than 1 object
  IF  $n[i] > 1$  THEN add  $\langle P, \bar{x} \rangle$  to the output buffer
reduce()
FOR each output buffer entry  $\langle P_i, \bar{x}_i \rangle$  DO output  $\langle 1/P_i, \bar{x}_i \rangle$ 

reduce() IS
FOR each output buffer entry  $\langle P_i, \bar{x}_i \rangle$  DO
  Let  $P_i' = \min \{ M / [\alpha_D \cdot n[i]^e \cdot \Gamma^\delta], 1 \}$ 
  WITH prob.  $P_i'/P_i$ , replace this entry with  $\langle P_i', \bar{x}_i \rangle$ 
  OTHERWISE remove this entry

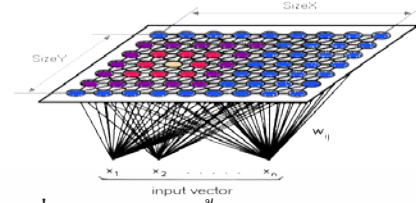
discordancy(lsd, n) IS
IF  $n \leq 1$  THEN  $n = 2$ 
IF  $\Gamma = \Gamma$  THEN
   $\text{disc} = \text{lsd} / ((n * (n-1)) / 2)$ 
ELSE
   $\text{disc} = \text{lsd} / (n-1)$ 
RETURN disc
    
```

อัลกอริทึม DBSPACE จะทำการแบ่งกลุ่มข้อมูล และสุ่มเลือกข้อมูลไปพร้อมๆ กันแบบคู่ขนานภายในการอ่านข้อมูลเพียงหนึ่งรอบ และข้อมูลที่ถูกเลือกจะถูกเก็บไว้ในบัฟเฟอร์ (output buffer) โดยในบัฟเฟอร์ประกอบด้วย $\langle P_i, \bar{x}_i \rangle$ ซึ่งแสดงถึงการที่ \bar{x}_i ถูกเลือกด้วยความน่าจะเป็นเท่ากับ P_i และมีการรับข้อมูลเข้ามาเรื่อยๆ \bar{x}_i จะมีความน่าจะเป็นที่จะถูกเลือกเท่ากับ P_i' ดังนั้นจึงต้องมีการปรับปรุงค่าในบัฟเฟอร์อยู่เสมอเพื่อให้ชุดข้อมูลสุ่มที่ได้เป็นปัจจุบัน โดยการเก็บ $\langle P_i', \bar{x}_i \rangle$ แทน $\langle P_i, \bar{x}_i \rangle$ ด้วยความน่าจะเป็น P_i' / P_i (หรือทำการลบอ็อบเจกต์นั้นออกจากบัฟเฟอร์)

2.2 Self-Organizing Map

Self-Organizing Map (SOM) เป็นนิเวศน์เน็ตเวิร์กแบบไม่มีผู้สอน [7] นำเสนอโดยศาสตราจารย์โคโฮเนน (Kohonen) ในปี ค.ศ. 1982 ประกอบด้วยเซลล์ 2 ชั้น ดังรูปที่ 3 ชั้นแรกคือชั้นของอินพุต (Input Layer) ประกอบด้วยเซตของอินพุตเวกเตอร์ $x(t)$ (1 x n มิติ) โดยที่ i คืออิน

เด็กซ์ของอินพุต ชั้นที่สองคือชั้นของเอาต์พุตประกอบไปด้วยโหนดของนิเวศน์เน็ตเวิร์กที่เรียงตัวอยู่ในรูปแบบของแผนภาพ 2 มิติ ในแต่ละโหนด i จะเป็นค่าเวกเตอร์น้ำหนักแทนด้วย $w_i(t)$ นั่นคือ $w_i(t) \in \mathcal{R}^n$ โดยที่ \mathcal{R}^n



รูปที่ 3 แสดงโมเดลพื้นฐานของ SOM

กระบวนการเรียนรู้ของ SOM เกิดขึ้นจากการปรับตัวของเวกเตอร์น้ำหนักที่มีต่ออินพุตเวกเตอร์ โดยเริ่มแรกจะกำหนดน้ำหนักเริ่มต้นขนาดเล็กให้กับทุกโหนด จากนั้นจะเริ่มต้นกระบวนการเรียนรู้ดังนี้

1. เลือกอินพุตเวกเตอร์จากอินพุตโดเมน
2. เปรียบเทียบอินพุตเวกเตอร์ $x(t)$ กับโหนด $w(t)$ ทุกโหนดเพื่อหาโหนดชนะจากโหนดทั้งหมด
3. ปรับเวกเตอร์น้ำหนักของโหนดชนะ เพื่อให้โหนดชนะเข้าใกล้อินพุตมากขึ้น
4. ปรับเวกเตอร์น้ำหนักของโหนดใกล้เคียง เพื่อให้อินพุตเวกเตอร์ถัดไปที่มีค่าใกล้เคียงกับ โหนดชนะ กระบวนการเหล่านี้จะถูกทำซ้ำไปเรื่อย ๆ จนกว่าจะสอดคล้องตามเงื่อนไข

กระบวนการเรียนรู้ข้างต้นมีการคำนวณที่สำคัญอยู่ 2 ส่วนคือ ส่วนแรกคือการคำนวณเพื่อหาโหนดชนะ c คือโหนดที่มีระยะห่างระหว่างอินพุตเวกเตอร์กับเวกเตอร์น้ำหนักน้อยที่สุด โดยใช้ฟังก์ชันวัดระยะทางแบบยูคลิด (Euclidean distance) โหนดชนะหาได้ดังสมการที่ 12

$$c : w.t = \min \| x(t) - w_i(t) \| \tag{12}$$

ส่วนที่สองคือการปรับเวกเตอร์น้ำหนักเพื่อให้เข้าใกล้อินพุตมากขึ้น การปรับค่าน้ำหนักแสดงดังสมการที่ 13

$$w_i(t+1) = w_i(t) + \alpha(t) \times h_c(t) \times [x(t) - w_i(t)] \tag{13}$$

เมื่อ t คือรอบปัจจุบันของการเรียนรู้, $x(t)$ คืออินพุตเวกเตอร์ปัจจุบัน, $w_i(t)$ คือเวกเตอร์น้ำหนัก, $\alpha(t)$ คืออัตราการเรียนรู้ โดยที่อัตราการเรียนรู้ $\alpha(t)$ จะขึ้นอยู่กับจำนวนรอบซึ่งแสดงเป็นสมการเชิงเส้น ดังสมการที่ 14

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$$\alpha(t) = \alpha(0) \times \frac{T-t}{T} \quad (14)$$

เมื่อ T คือจำนวนรอบทั้งหมด, t คือจำนวนรอบปัจจุบัน h_{ci} คือฟังก์ชันที่ใช้ในการกำหนดน้ำหนักในการปรับค่าโหนดใกล้เคียงโดยทั่วไปแล้ว จะใช้ฟังก์ชันเกาส์เซียน ซึ่งสามารถเขียนได้ดังสมการที่ 15

$$h_{ci}(t) = \exp\left(-\frac{\|r_c - r_i\|^2}{2\sigma^2(t)}\right) \quad (15)$$

โดยปกติรัศมีของโหนดใกล้เคียงจะค่อย ๆ ลดลงตามจำนวนรอบในการเรียนรู้ t ดังสมการที่ 16

$$\sigma(t+1) = 1 + (\sigma(t) - 1) \times \frac{T-t}{T} \quad (16)$$

2.3 ขั้นตอนของวิธีที่นำเสนอ

ข้อมูลอินพุตที่ใช้เรียนรู้ในแผนภาพมี 2 ส่วนโดยส่วนแรกเป็นข้อมูลเอาต์พุตที่ได้จาก DBSPACE และส่วนที่ 2 เป็นข้อมูลอินพุตทั้งหมดโดยข้อมูลส่วนแรกเรียนรู้ประมาณ 70 เปอร์เซ็นต์ของจำนวนรอบการเรียนรู้สูงสุด และข้อมูลส่วนที่ 2 จะเรียนรู้อีกประมาณ 30 เปอร์เซ็นต์ซึ่งขั้นตอนการทำงานแสดงได้ดังนี้

กำหนดให้ชุดข้อมูลฝึกฝนอินพุตเป็น x_1, x_2, \dots, x_K เมื่อ $x_k = [x_{k1}, \dots, x_{kr}]^T$, $k=1, \dots, K$, และ K เป็นจำนวนข้อมูลทั้งหมดในชุดฝึกฝน L_{max} เป็นจำนวนรอบสูงสุดในการเรียนรู้ที่กำหนด

ขั้นตอนที่ 1. นำข้อมูลอินพุต x_1, x_2, \dots, x_K ผ่านกรรมวิธีการลดข้อมูลโดย DBSPACE ได้ข้อมูลสุ่มชุดหนึ่ง y_1, y_2, \dots, y_J ซึ่งเป็นเอาต์พุตของ DBSPACE

ขั้นตอนที่ 2. สุ่มค่าน้ำหนัก w_i , $i=1, 2, \dots, M$

ขั้นตอนที่ 3. รับข้อมูลชุดฝึกสอนจากเอาต์พุตของ DBSPACE เข้าโครงข่ายคำนวณระยะทางระหว่าง $[y_j - w_i]$

ขั้นตอนที่ 4. หาค่านิวรอนที่ชนะจากนั้นปรับค่าเวกเตอร์น้ำหนักของประสาทเทียม

ขั้นตอนที่ 5. ถ้า $j < J$ ให้ $j=j+1$ ทำการคำนวณซ้ำจากขั้นตอนที่ 3

ขั้นตอนที่ 6. ถ้าจำนวนรอบ $m < (L_{max} - L_{max} * 0.3)$ (0.3 คือ 30 เปอร์เซ็นต์) กำหนด $j=1, m=m+1$ เริ่มรอบฝึกฝนใหม่จากขั้นตอนที่ 3

ขั้นตอนที่ 7. รับข้อมูลชุดฝึกสอน x_1, x_2, \dots, x_K เข้า

โครงข่ายคำนวณระยะทางระหว่าง $[x_k - w_i]$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ขั้นตอนที่ 8. หาค่าประสาทเทียมที่ชนะปรับค่าเวกเตอร์น้ำหนักของประสาทเทียม

ขั้นตอนที่ 9. ถ้า $k < K$ ให้ $k=k+1$ ทำการคำนวณซ้ำจากขั้นตอนที่ 7

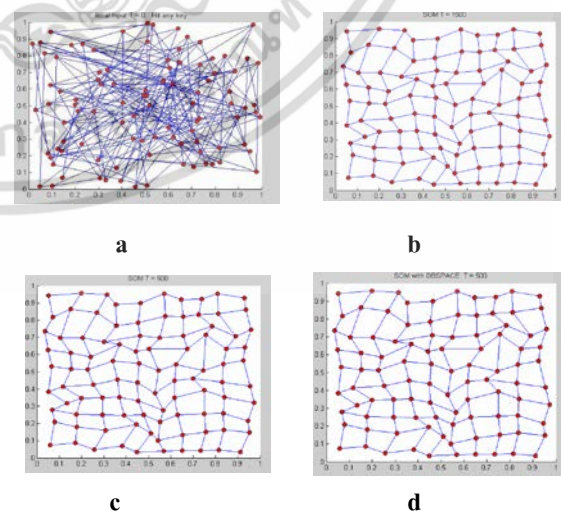
ขั้นตอนที่ 10. ถ้าจำนวนรอบ $m < L_{max}$ กำหนด $k=1, m=m+1$ เริ่มรอบฝึกฝนใหม่จากขั้นตอนที่ 7

3. การทดลองและผลการทดลอง

ทำการทดลองโดยเครื่องคอมพิวเตอร์มีหน่วยประมวลผลเป็น Intel Core2Duo 2.0GHz Memory 2 GB ระบบปฏิบัติการ Windows XP

3.1 การทดลองที่ 1

ชุดข้อมูลที่ใช้มีค่าระหว่าง 0-1 โดยการสุ่ม เป็นเวกเตอร์ 2 มิติจำนวน 1000 เวกเตอร์ นำมาทำการเรียนรู้ด้วย SOM แบบเดิมและแบบที่นำเสนอขนาดโมเดล 10×10 โดยกำหนดอัตราการเรียนรู้เริ่มต้น $\alpha = 0.5$ จะเห็นได้ว่าเมื่อทำการเรียนรู้ในแผนภาพ SOM แบบเดิมจำนวน 1500 รอบรูปที่ 4b และ 500 รอบรูปที่ 4c แผนภาพของโหนดที่แสดงไม่แตกต่างกันแสดงว่าการเรียนรู้เพียง 500 รอบก็เพียงพอในการทดสอบและเมื่อนำข้อมูลมาเรียนรู้ผ่าน SOM ที่นำเสนอโดยเรียนรู้ 500 รอบ แผนภาพที่ได้เหมือนกันกับที่เรียนรู้โดย SOM แบบเดิมดังรูปที่ 4d ส่วนเวลาที่ใช้ในการเรียนรู้ลดลงแสดงดังตารางที่ 2



รูปที่ 4 แสดงแผนภาพของ SOM

ตารางที่ 2 ผลการทดลองในการเรียนรู้ข้อมูลส้ม

วิธีการ	จำนวนรอบ	เวลา(วินาที)
SOM	500	150.89
Our Method	500	67.62

3.2 การทดลองที่ 2

ชุดข้อมูล ecoli[8] ซึ่งมีข้อมูลทั้งหมด 336 ชุด 8 แอทธิบิตโดยนำมาทดลองจำนวน 5 กลุ่ม จำนวน 327 ข้อมูล ทดลองโดยผ่าน SOM แบบเดิมและแบบที่นำเสนอ ขนาดโมเดล 4x4 โดยกำหนดอัตราการเรียนรู้เริ่มต้น $\alpha = 0.5$ และวัดประสิทธิภาพโดยการวัดค่าแบบเอนโทรปี (entropy) คือ การวัดโดยอาศัยความน่าจะเป็นของข้อมูลที่อยู่ในกลุ่ม สามารถหาได้จาก

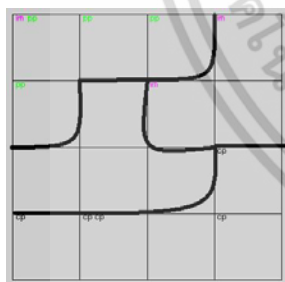
$$E_j = -\sum_i p_{ij} \log(p_{ij}) \tag{16}$$

โดยที่ค่าเอนโทรปีรวมของทุกกลุ่มสามารถหาได้จาก

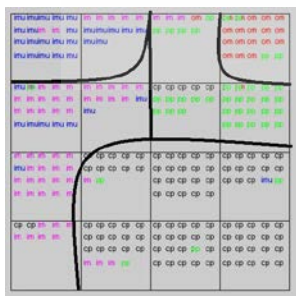
$$E_{cs} = \sum_{j=1}^m \frac{n_j \times E_j}{n} \tag{17}$$

โดยที่ p คือความน่าจะเป็นที่สมาชิกของโหนด ขึ้นอยู่กับกลุ่ม j, n_j คือ จำนวนข้อมูลทั้งหมดในโหนด j, n คือ จำนวนเอกสารทั้งหมด

ได้ผลการทดลองดังตารางที่ 3 และแผนภาพดังรูปที่ 5



a) แสดงแผนภาพหลังจากเรียนรู้ด้วยอินพุทจากDBSPACE



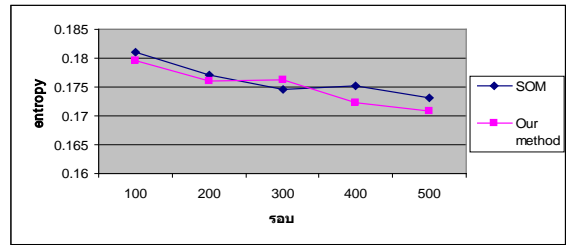
b) Som



c) Our method

รูปที่ 5 แสดงแผนภาพที่ได้จากข้อมูลชุด ecoli

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการวิจัย ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 6 กราฟของจำนวนรอบการทำงานและค่าเอนโทรปี

ตารางที่ 3 ผลการทดลองในการเรียนรู้ชุดข้อมูล ecoli

วิธีการ	จำนวนรอบ	เวลา(วินาที)	ค่าเอนโทรปี
SOM	1500	42.062	0.17149
Our Method	1500	16.156	0.16702

4.สรุป

การทดลองแสดงว่าวิธีที่นำเสนอใช้เวลาในการเรียนรู้ลดลงเมื่อเปรียบเทียบกับ SOM แบบเดิมเนื่องจากการประยุกต์ใช้ DBSPACE นั้นช่วยให้จำนวนข้อมูลที่ใช้เรียนรู้ลดลงและข้อมูลที่ได้เป็นตัวแทนกลุ่มของข้อมูลทั้งหมด เมื่อนำข้อมูลนั้นไปฝึกในโครงข่ายก่อนใช้ข้อมูลเดิม จึงส่งผลให้การจัดเรียงตัวของแผนภาพทำได้เร็วขึ้น และใช้เวลาในการเรียนรู้ลดลงมากกว่า 50 เปอร์เซ็นต์แต่ประสิทธิภาพในการจัดกลุ่มข้อมูลยังเหมือนเดิม

5.เอกสารอ้างอิง

- [1] J.Kangas. "Time dependent self-organizing maps for Speech recognition." ICANN-91, Espoo, June 24-28, 1991
- [2] Kaski S., Honkela T., Lagus K., and Kohonen T. "WEBSOM Self-organizing maps of document collections" Neurocomputing, Volume 21, 1998, pp. 101-117.
- [3] Samuel Kaski, Jari Kangas, Teuvo Kohonen. "Bibliography of Self-Organizing Map(SOM) Paper: 1981-1977" Neural Computing Survey I, 1998.
- [4] Mu-Chunn Su, Hsiao-Te Chang. "Fast Self-Organizing Feature Map Algorithm." IEEE Transaction on neural Networks vol.11, no.3, May 2000.
- [5] Kaski, S. "Fast winner search for SOM-based monitoring and retrieval of high-dimensional data" Proceedings of ICANN99 Ninth International Conference on Artificial Neural Networks, London, vol 2, pp940-945
- [6] ธรรมศักดิ์ เรือรณีเวศน์, "การลดขนาดข้อมูลด้วยน้ำหนักความหนาแน่นเพื่อการจัดกลุ่มของมูลขนาดใหญ่" วิทยานิพนธ์ปริญญาโทฉบับที่ ๓๓ มหาวิทยาลัยสุรนารี 2548
- [7] Helge Ritter, Thomas Martinetz and Klaus Schulten. Neural computation and organizing maps: an introduction Massachusetts: Addison-wesley. 1992
- [8] <http://archive.ics.uci.edu/ml/datasets/ecoli> ด้านการค้า