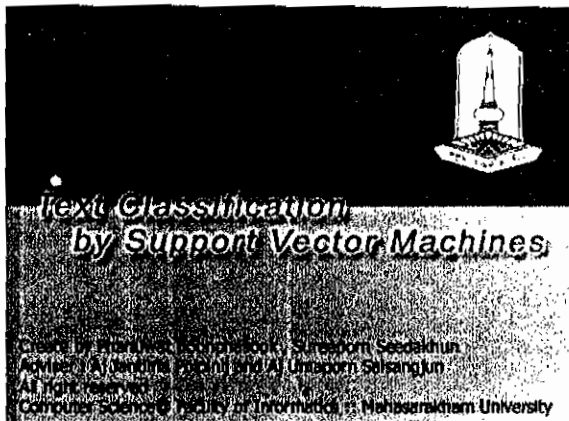


การจัดกลุ่มเอกสารข้อความภาษาไทยด้วยอัลกอริทึม SVMs

(Thai Text Classification By Support Vector Machines)



ได้รับทุนอุดหนุนจากโครงการวิจัย พัฒนา และ วิศวกรรม

โครงการแข่งขันพัฒนาโปรแกรมคอมพิวเตอร์

แห่งประเทศไทย ครั้งที่ 7 ประจำปีงบประมาณ 2547

บทคัดย่อ

เนื่องจากความนิยมในการใช้งานระบบอินเทอร์เน็ตและ ความหลากหลายที่เกี่ยวกับบริการบนระบบอินเทอร์เน็ตที่เพิ่มมากขึ้น ก่อให้เกิดปัญหาตามมาในเรื่องของจำนวนสารสนเทศที่เพิ่มมากขึ้น การจัดกลุ่มเอกสารเป็นอีกหนึ่ทางเลือกที่สามารถช่วยให้ผู้ใช้ระบบ สามารถค้นหาสารสนเทศที่ต้องการได้รวดเร็วมากขึ้น สำหรับงานวิจัย ฉบับนี้ได้นำเสนอวิธีการจัดกลุ่มเอกสารด้วยอัลกอริทึม Support Vector Machines (SVMs) สำหรับการสร้างตัวแยกแยะเอกสารเพื่อจัด กลุ่มเอกสารข้อความภาษาไทยแบบอัตโนมัติ บนพื้นฐานของข้อมูลที่มี ใช้ช่่องริง

หลักการและเหตุผล

เนื่องจากการเพิ่มขึ้นของจำนวนเอกสารซึ่งมีอยู่เป็น จำนวนมาก ได้เริ่มสร้างปัญหาให้กับผู้สืบค้นข้อมูลสารสนเทศเพื่อ นำไปใช้งาน ดังนั้นงานทางด้านการค้นหาเอกสารแบบอัตโนมัติซึ่ง ได้รับความสะดวกแก่นักวิจัยที่ทำงานในด้านการค้นหาสารสนเทศ (Information Retrieval) ซึ่งหนึ่งในงานวิจัยที่สำคัญ คือ การจำแนก ประเภทของเอกสาร (Text Classification)

ในการจัดกลุ่มเอกสารกับข้อความภาษาไทยจะค่อนข้าง ยุ่งยากกว่าเอกสารที่เป็นภาษาอังกฤษ (ซูริวัตน์และคณะ, 2000) เนื่องจากภาษาไทยมีลักษณะภาษาที่มีการเขียนในรูปแบบสายอักขระ ก่อเนื่องที่ไม่มีจุดสิ้นสุดของประโยคอย่างชัดเจน ทำให้ปัญหาในเรื่อง การสร้างคำขอ (Query) ที่เป็นส่วนสำคัญในการจัดกลุ่มเอกสาร

เนื่องจากการเพิ่มขึ้นของจำนวนเอกสารซึ่งมีอยู่เป็น จำนวนมาก ได้เริ่มสร้างปัญหาให้กับผู้สืบค้นข้อมูลสารสนเทศเพื่อ นำไปใช้งาน ดังนั้นงานทางด้านการค้นหาเอกสารแบบอัตโนมัติซึ่ง ได้รับความสะดวกแก่นักวิจัยที่ทำงานในด้านการค้นหาสารสนเทศ (Information Retrieval) ซึ่งหนึ่งในงานวิจัยที่สำคัญ คือ การจำแนก ประเภทของเอกสาร (Text Classification)

ในการจัดกลุ่มเอกสารกับข้อความภาษาไทยจะค่อนข้างยุ่งยาก กว่าเอกสารที่เป็นภาษาอังกฤษ (ซูริวัตน์และคณะ, 2000) เนื่องจาก ภาษาไทยมีลักษณะภาษาที่มีการเขียนในรูปแบบสายอักขระต่อเนื่องที่ไม่มี จุดสิ้นสุดของประโยคอย่างชัดเจน ทำให้ปัญหาในเรื่องของการสร้างคำขอ (Query) ที่เป็นส่วนสำคัญในการจัดกลุ่มเอกสาร

สำหรับโครงงานนี้มีเหตุผลอีกกรณีคือ Support Vector Machines (SVMs) มีประสิทธิภาพสูงในการจัดกลุ่มเอกสารข้อความภาษาไทย เพื่อเป็นอีกแนวทางหนึ่งในการจัดกลุ่มเอกสารโดยเฉพาะเอกสารข้อความ ภาษาไทย

วัตถุประสงค์และเป้าหมายของการจัดทำโครงการ

เพื่อพัฒนาระบบการจัดกลุ่มเอกสารข้อความภาษาไทย (Thai Text Classification) แบบอัตโนมัติ โดยการใช้วิธีการแมชชีนเลิร์นนิง Machine Learning ที่เรียกว่า Support Vector Machines (SVMs) และเพื่อ เพิ่มประสิทธิภาพในการค้นหาสารสนเทศ

ขอบเขตและข้อจำกัดของโครงการ

พัฒนาระบบการจัดกลุ่มเอกสารข้อความภาษาไทย (Thai Text Classification) ด้วยอัลกอริทึมแมชชีน Machine Learning ที่เรียกว่า Support Vector Machines (SVMs) โดยเอกสารที่นำมาใช้ในการทดลองการจัดกลุ่ม เอกสารจะเป็นเอกสารข้อความภาษาไทยเท่านั้น ซึ่งลักษณะของเอกสารจะ จัดให้อยู่ในรูปแบบ HTML ที่ไม่รวมถึงภาพ เสียงลักษณะพิเศษ สูตรหรือสมการ ใดๆ

สมการหลัก

$$wx + b$$

ถ้าค่าของ $wx + b > 0$ จะกำหนดให้ค่า $y = 1$ ซึ่งจะจัดอยู่ใน Class 1
ถ้าค่าของ $wx + b < 0$ จะกำหนดให้ค่า $y = -1$ ซึ่งจะจัดอยู่ใน Class 2

การวัดประสิทธิภาพ

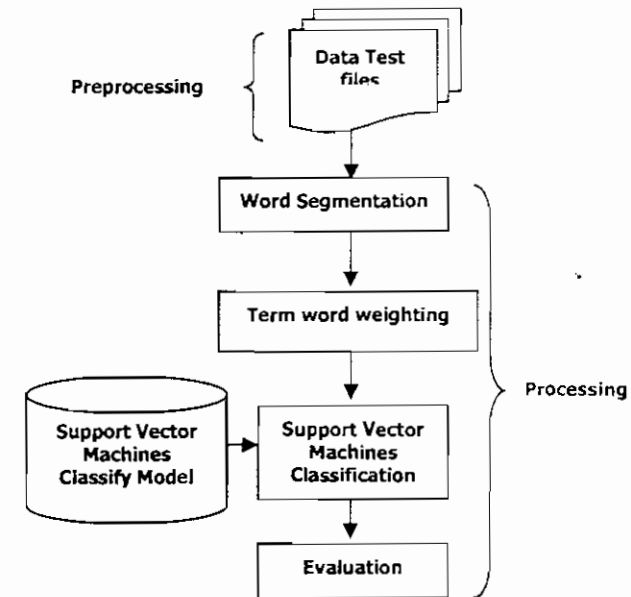
$$\text{ค่าความแม่นยำ} = \frac{\text{จำนวนเอกสารที่ถูกจัดกลุ่มที่ค้นคืนได้}}{\text{จำนวนเอกสารที่ถูกจัดกลุ่มทั้งหมดที่ค้นคืนออกมาได้}}$$

$$\text{ค่าความระลึก} = \frac{\text{จำนวนเอกสารที่ถูกจัดกลุ่มที่ค้นคืนได้}}{\text{จำนวนเอกสารที่ถูกจัดกลุ่มทั้งหมดในฐานข้อมูล}}$$

การวัดคุณภาพซึ่งเป็นค่าผสมระหว่างค่าความแม่นยำและค่าความ ระลึก

$$F_{1,2} = \frac{(b^2 + 1)P_{1,1}R_{1,1}}{b^2P_{1,1} + R_{1,1}}$$

ขั้นตอนการดำเนินการ

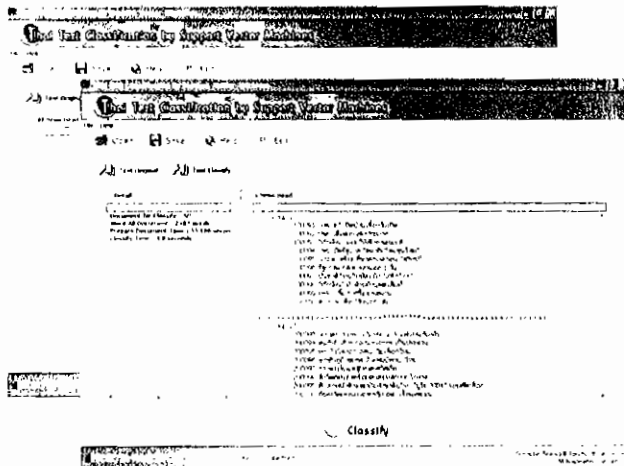


สรุปผลการทดลอง

ผลการประเมินโปรแกรมที่ได้พัฒนาขึ้น โดยได้นำไปทดสอบจัดกลุ่มเอกสารข่าว จำนวน 100 เอกสาร โดยแบ่งออกเป็น 5 ประเภท คือ ข่าวการเมือง, ข่าวกีฬา, ข่าวเศรษฐกิจ, ข่าวบันเทิง และข่าวต่างประเทศ หลังจากนั้นจะเป็นการประเมินผลความถูกต้องของผลการจัดกลุ่มเอกสารโดยโปรแกรมที่พัฒนาขึ้นดังกล่าว สามารถสรุปได้ว่า ระบบมีความน่าเชื่อถือในระดับที่น่าพอใจคือประมาณ 72.0%

ตัวอย่างโปรแกรม

ตัวอย่างโปรแกรมที่ทำการจัดกลุ่มเอกสาร



ผู้จัดทำ

นายภาณุวัฒน์ บุญมาสุข
(phanuwat_cs@ hotmail.com)
นางสาวสุรีย์พร สีดาคุณ
(Victoria_cs99@ hotmail.com)

อาจารย์ที่ปรึกษา

อุมารณีย์ สายแสงจันทร์
(umaporn.sa@msu.ac.th)
จันทิมา พลพินิจ
(jantima.p@msu.ac.th)

สาขาวิชาวิทยาการคอมพิวเตอร์ คณะวิทยาการสารสนเทศ

มหาวิทยาลัยมหาสารคาม

อำเภอกันทรวิชัย จังหวัดมหาสารคาม 44150

โทรศัพท์ :: 0-4375-4359 หรือ

0-4375-4322-40 ต่อ 2414 , 2453

แฟกซ์ :: 0-4375-4359