

Computer Science Department
Faculty of Informatics, Maharakham University

โปรแกรมโครงงาน

A Method of Imbalanced Sentiment Classification

กระบวนการสำหรับการจำแนกความรู้สึกที่มีข้อมูลไม่สมดุล



MAHASARAKHAM UNIVERSITY



ผู้พัฒนา : พีระวัฒน์ บุญบ้านจิว (Pheerawat Bunbannjio)

อาจารย์ที่ปรึกษา : จันทิมา พลพิณ (Jantima Polpinij)

Intellect Laboratory สาขาวิทยาการคอมพิวเตอร์ คณะวิทยาการสารสนเทศ มหาวิทยาลัยมหาสารคาม
knarf.pheerawat@gmail.com, jantima.p@msu.ac.th

ที่มาและความสำคัญ

การจำแนกความรู้สึก (Sentiment Classification) คือการจำแนกเอกสารตามข้อความซึ่งโดยทั่วไปจะจำแนกเป็นความรู้สึกที่เป็นบวก (Positive) ความรู้สึกที่เป็นลบ (Negative) และความรู้สึกที่เป็นกลาง (Neutral) โดยการจำแนกความรู้สึกนั้น ได้รับการศึกษาอย่างต่อเนื่อง เพราะการประยุกต์ใช้ในหลายลักษณะ แต่โดยทั่วไปมักจะนิยมใช้ในการจำแนกความรู้สึกที่มีการแสดงไว้ในรูปแบบข้อความ (Text) เช่น ประยุกต์ใช้ในการจัดอันดับความรู้สึกจากข้อความแสดงความคิดเห็นของผู้คนที่ติดต่อสินค้าและบริการ การประยุกต์ใช้เพื่อวิเคราะห์ความรู้สึกของผู้คนในเรื่องการเมือง เป็นต้น ซึ่งปัญหาความไม่สมดุลของข้อมูลในคลาสนั้น เกิดจากกลุ่มตัวอย่างที่ใช้ในการเรียนรู้มีข้อมูลไม่สมดุลกัน โดยกลุ่มที่มีข้อมูลมากกว่าจะเรียกว่า "ข้อมูลกลุ่มหลัก (Majority Classes)" ขณะที่กลุ่มตัวอย่างที่มีข้อมูลจำนวนน้อยกว่าจะเรียกว่า "ข้อมูลกลุ่มรอง (Minority Class)" เนื้อหาของข้อมูลในลักษณะนี้ไปเรียนรู้เพื่อสร้างตัวจำแนกความรู้สึก (Sentiment Classifier) ข้อมูลใหม่ๆ ที่อ่านเข้ามาเพื่อวิเคราะห์เพื่อจำแนกกลุ่มด้วยตัวจำแนกความรู้สึกดังกล่าว ก็มีแนวโน้มที่จะทำมากลุ่มของข้อมูลในไปถึงทิศทางของข้อมูลกลุ่มหลักที่ใช้ในการเรียนรู้ตัวจำแนกความรู้สึก ดังนั้น ในโครงการปริญญาโทฉบับนี้มี จึงได้นำเสนอการศึกษาร่วมกันเพื่อหาความไม่สมดุลของข้อมูลในการจำแนกความรู้สึกด้วยเทคนิคการให้น้ำหนักค่า 5 เทคนิค คือ TF-IDF, Delta TF-IDF, TF-IDF-ICF, TF-RF และ TF-IGM ร่วมกับแบบซิมเพล็กซ์ 3 ตัว คือ Naive Bayes, K-Nearest Neighbor และสุดท้าย Convolution Neural Network

คำสำคัญ: การจำแนกเอกสาร, การให้น้ำหนักค่า, ข้อมูลไม่สมดุล, ซัพพอร์ตเวกเตอร์แมชชีน

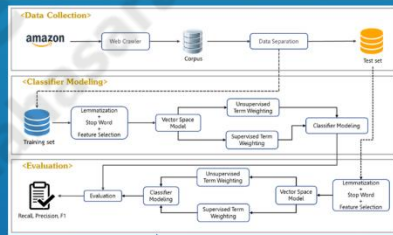
วัตถุประสงค์

นำเสนอกระบวนการสำหรับการจำแนกความรู้สึกที่มีข้อมูลไม่สมดุลโดยมีเครื่องมือหลักเทคนิคการให้น้ำหนักค่าแบบมีผู้สอน

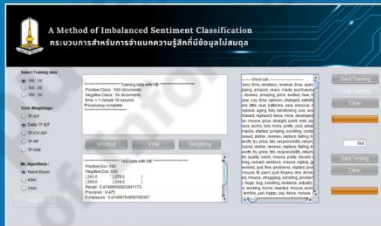
ประโยชน์ที่คาดว่าจะได้รับ

ได้กระบวนการในการจำแนกข้อความแสดงความรู้สึกที่มีข้อมูลแบบไม่สมดุล

กรอบการดำเนินงาน



ตัวอย่างหน้าจอการทำงาน



สรุป

เนื่องจากบ่อยครั้งที่ การจำแนกเอกสารที่ไม่สมดุลกันนั้นมีการเอนเอียงการให้ค่าแบบไม่เชิงที่ข้อมูลมากกว่าเนื่องจากมีข้อมูลที่ครอบคลุมการทำนายที่ดีกว่า ดังนั้นงานวิจัยฉบับนี้จึงได้นำเสนอวิธีการการจำแนกข้อมูลที่ไม่สมดุลด้วยการให้น้ำหนักค่าเปรียบเทียบ 2 รูปแบบหลักคือ UTW และ STW โดย UTW ใช้รูปแบบการให้น้ำหนักค่าที่ได้รับค่านิยมมากที่สุดคือ TF-IDF และ STW ใช้ทั้งหมด 4 รูปแบบคือ Delta TF-IDF, TF-ICF-IDF, TF-RF และ TF-IGM โดยผลที่ได้คือการให้น้ำหนักค่าแบบ STW มีประสิทธิภาพในการจำแนกข้อมูลที่ไม่สมดุลมากกว่ารูปแบบการให้น้ำหนักค่าแบบ UTW ซึ่งได้แก่การให้น้ำหนักค่าแบบ TF-IGM โดยใช้อัลกอริทึม CNN ในการสร้างโมเดล มีค่าเฉลี่ย F-measure สูงที่สุดอยู่ที่ 74.41%