

การลดพีเจอร์โดยใช้ต้นไม้วลีสำหรับการจำแนกประเภทเอกสาร

PHRASE-TREE-BASED FEATURE REDUCTION FOR TEXT
CATEGORIZATION

ปรีสุทธิ์ จิตต์ภักดี
PARISUT JITPAKDEE

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาค้นคว้าตามหลักสูตรปริญญาวิทยาศาสตรมหาบัณฑิต

สาขาวิชาเทคโนโลยีสารสนเทศ

บัณฑิตวิทยาลัย

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

พ.ศ. 2550

KMITL-2007-IT-M-001-123

สำนักหอสมุดกลาง พระจอมเกล้าลาดกระบัง

การลดพีเจอร์โดยใช้ต้นไม้วิไลสำหรับการจำแนกประเภทเอกสาร

**PHRASE-TREE-BASED FEATURE REDUCTION FOR TEXT
CATEGORIZATION**

ปรีสุทธิ์ จิตต์ภักดี

PARISUT JITPAKDEE

เลขหมู่.....
เลขทะเบียน.....**77950**
วัน,เดือน,ปี.....**1.2.ค.พ.2551**

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรมหาบัณฑิต

สาขาวิชาเทคโนโลยีสารสนเทศ

บัณฑิตวิทยาลัย

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

พ.ศ. 2550

KMITL-2007-IT-001-128

**PHRASE-TREE-BASED FEATURE REDUCTION FOR TEXT
CATEGORIZATION**

PARISUT JITPAKDEE

**A THESIS SUBMITTED IN PARTIAL FULFILLMENT
OF THE REQUIREMENT FOR THE DEGREE OF
MASTER OF SCIENCE IN INFORMATION TECHNOLOGY
SCHOOL OF GRADUATE STUDIES
KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG
2007
KMITL-2007-IT-001-128**

COPYRIGHT 2007

SCHOOL OF GRADUATE STUDIES

KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

หัวข้อวิทยานิพนธ์	การลดพีเจอร์โดยใช้ต้นไม้วลิสำหรับการจำแนกประเภทเอกสาร
นักศึกษา	นางสาวปรีสุทธิ์ จิตต์ภักดิ์
รหัสประจำตัว	47066405
ปริญญา	วิทยาศาสตรมหาบัณฑิต
สาขาวิชา	เทคโนโลยีสารสนเทศ
พ.ศ.	2550
อาจารย์ที่ปรึกษาวิทยานิพนธ์	รศ.ดร.วรพจน์ กรีสู่ระเดช

บทคัดย่อ

วิทยานิพนธ์ฉบับนี้ นำเสนอวิธีการและอัลกอริทึมใหม่ในการลดพีเจอร์โดยใช้ต้นไม้วลิสำหรับการจำแนกประเภทเอกสาร ในการจำแนกประเภทเอกสาร เอกสารที่เรารู้ประเภทของพวกมันแล้ว จะถูกแทนให้อยู่ในรูปของเวกเตอร์ของพีเจอร์ พีเจอร์เหล่านี้เป็นกุญแจสำคัญของประสิทธิภาพในการจำแนกประเภทเอกสาร เซตของพีเจอร์ที่ดีจะต้องให้ความถูกต้องในการจำแนกประเภทสูง และมีขนาดเล็ก โดยการใช่วลิเป็นพีเจอร์นั้นดีกว่าการใช้คำเดี่ยว เพราะวลิสามารถให้ความหมายได้มากกว่าคำเดี่ยว ซึ่งต้นไม้วลิที่นำเสนอนี้จะสร้างพีเจอร์วลิ และแสดงถึงความสัมพันธ์ระหว่างวลิและวลิย่อยของพวกมัน ซึ่งช่วยลดความซ้ำซ้อนของพีเจอร์ทำให้ได้พีเจอร์ที่มีประสิทธิภาพ จากการทดลองเบื้องต้น วิธีการที่นำเสนอนี้ให้เซตของพีเจอร์ที่มีขนาดเล็กกว่าวิธีการแบบเดิม และยังทำให้ความถูกต้องในการจำแนกประเภทเอกสารเพิ่มขึ้นอีกด้วย

Thesis Title	Phrase-Tree-based Feature Reduction for Text Categorization
Student	Ms. Parisut Jitpakdee
Student ID.	47066405
Degree	Master of Science
Programme	Information Technology
Year	2007
Thesis Advisor	Assoc.Prof.Dr.Worapoj Kreesuradej

ABSTRACT

This thesis proposes a new method and an algorithm for feature reduction based on phrase tree for text categorization. In text categorization, the documents that we know their categories will be represented as vectors of features. These features are the key of the performance of text categorization. The good features set must give the high accuracy of categorization and have small size. Using phrases as features is better than using words, because phrase gives the meaning more than single word. This proposed phrases tree generates phrase features and represents the relationship of these phrases and their sub phrases that help to remove redundant of feature to make the efficient features set. From preliminary experiments, the proposed method gives smaller set of features than conventional methods. Furthermore, the classification model obtained from the propose method provides higher accuracy.

กิตติกรรมประกาศ

วิทยานิพนธ์ฉบับนี้สำเร็จลุล่วงได้อย่างดี ด้วยการให้คำปรึกษาพร้อมทั้งคำชี้แนะแนวทางการแก้ปัญหาของงานวิจัยตั้งแต่เริ่มต้นจนเสร็จสมบูรณ์จาก รศ.ดร.วรพจน์ กรีสระเดช ซึ่งเป็นอาจารย์ผู้ควบคุมวิทยานิพนธ์ ผู้วิจัยรู้สึกซาบซึ้งในความอนุเคราะห์ของท่านและขอขอบพระคุณเป็นอย่างสูง

ขอขอบพระคุณคุณย่าและญาติพี่น้องทุกท่าน ที่คอยสนับสนุนและคอยให้กำลังใจมาโดยตลอด

ขอขอบคุณคณาจารย์คณะเทคโนโลยีสารสนเทศทุกท่าน ที่ได้ประสิทธิ์ประสาทวิชาความรู้ต่างๆ เพื่อนำไปใช้ประโยชน์ในการทำวิจัย

ขอขอบคุณเจ้าหน้าที่คณะเทคโนโลยีสารสนเทศทุกท่าน ที่ให้คอยให้ความสะดวกในการทำงาน

ขอขอบคุณพี่ๆ และเพื่อนๆ ห้องปฏิบัติการ Data Mining & Data Exploration Laboratory ที่คอยให้ความช่วยเหลือ คำแนะนำ และให้กำลังใจซึ่งกันและกันตลอดมา

อนึ่ง งานวิจัยที่นำเสนอในงานวิทยานิพนธ์ฉบับนี้นั้น ส่วนหนึ่งกระทำภายใต้ห้องปฏิบัติการ Data Mining & Data Exploration Laboratory คณะเทคโนโลยีสารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

คุณค่าและประโยชน์จากวิทยานิพนธ์นี้ ผู้วิจัยขอมอบแด่ผู้มีพระคุณทุกท่าน

ปรีสุทธิ์ จิตต์ภักดี

สารบัญ

	หน้า
บทคัดย่อภาษาไทย	I
บทคัดย่อภาษาอังกฤษ	II
กิตติกรรมประกาศ.....	III
สารบัญ	IV
สารบัญตาราง	IX
สารบัญรูป	XV
บทที่ 1 บทนำ	1
1.1 ความเป็นมาและความสำคัญของปัญหา	1
1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา	2
1.3 ทฤษฎีหรือแนวคิดที่ใช้ในการวิจัย	3
1.4 ขอบเขตของการวิจัย	3
1.5 ขั้นตอนการศึกษา.....	3
บทที่ 2 ทฤษฎีพื้นฐานและงานวิจัยที่เกี่ยวข้อง.....	5
2.1 การจำแนกประเภทเอกสาร (Text Categorization)	5
2.2 การเตรียมเอกสาร (Document Pre-processing)	7
2.3 การสกัดฟีเจอร์ (Feature Extraction).....	7
2.4 การเลือกฟีเจอร์ (Feature Selection).....	8
2.5 การแทนเอกสาร (Document Representation)	11
2.6 วิธีการจำแนกประเภทเอกสาร (Categorization Method).....	13
2.7 การจำแนกประเภทเอกสารโดยใช้วลี	15
บทที่ 3 การลดฟีเจอร์โดยใช้ต้นไม้วลีสำหรับการจำแนกประเภทเอกสาร	19
3.1 การเตรียมเอกสาร (Document Pre-processing)	19
3.2 การสกัดวลี (Phrase Extraction).....	20

สารบัญ (ต่อ)

	หน้า
3.3 การลดพีเจอร์โดยใช้ต้นไม้วลีสำหรับการจำแนกประเภทเอกสาร	20
3.3.1 การสร้างต้นไม้วลี (Phrase-tree Construction)	21
3.3.2 การลดพีเจอร์ระดับ Local.....	22
3.3.3 การลดพีเจอร์ระดับ Global.....	23
3.4 การแทนเอกสาร (Document Representation)	25
บทที่ 4 การทดลองและผลการทดลอง.....	27
4.1 การวัดความถูกต้องของการจำแนกประเภทเอกสาร	27
4.1.1 ตัววัดอัตราความถูกต้อง (Accuracy Rate).....	27
4.1.2 ตัววัดอัตราความถูกต้อง F-measure.....	27
4.2 การออกแบบการทดลอง	29
4.2.1 การสุ่มตัวอย่างเอกสาร (Document Sampling)	29
4.2.2 การตรวจสอบไขว้ (Cross-Validation).....	30
4.2.3 รูปแบบของพีเจอร์ที่ใช้ในการทดลอง	30
4.3 ชุดเอกสารที่ใช้ในการทดลอง	31
4.4 ชุดเอกสาร PDDP F-series ที่ใช้ในการทดลองและผลการทดลอง	31
4.4.1 ชุดเอกสาร PDDP F-series ที่ใช้ในการทดลอง.....	31
4.4.2 ผลการทดลองของชุดเอกสาร PDDP F-series เมื่อใช้พีเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลี.....	31
4.4.3 ผลการทดลองของชุดเอกสาร PDDP F-series เมื่อใช้พีเจอร์แบบคำเดี่ยว.....	33
4.4.4 ผลการทดลองของชุดเอกสาร PDDP F-series เมื่อใช้พีเจอร์แบบคำเดี่ยวและไป-แกรม	35
4.4.5 ผลการทดลองของชุดเอกสาร PDDP F-series เมื่อใช้พีเจอร์แบบคำเดี่ยวและวลี	36
4.4.6 เปรียบเทียบผลการทดลองเมื่อใช้พีเจอร์แบบต่างๆ ของชุดเอกสาร PDDP F-series.....	36

สารบัญ (ต่อ)

	หน้า
4.5 ชุดเอกสาร PDDP J-series ที่ใช้ในการทดลองและผลการทดลอง	42
4.5.1 ชุดเอกสาร PDDP J-series ที่ใช้ในการทดลอง	42
4.5.2 ผลการทดลองของชุดเอกสาร PDDP J-series เมื่อใช้พีเจเออร์ที่ได้ จากการเลือกโดยใช้ต้นไม้วิลิ.....	42
4.5.3 ผลการทดลองของชุดเอกสาร PDDP J-series เมื่อใช้พีเจเออร์แบบคำ เดี่ยว.....	43
4.5.4 ผลการทดลองของชุดเอกสาร PDDP J-series เมื่อใช้พีเจเออร์แบบคำ เดี่ยวและ ไป-แกรม.....	46
4.5.5 ผลการทดลองของชุดเอกสาร PDDP J-series เมื่อใช้พีเจเออร์แบบคำ เดี่ยวและวลี	47
4.5.6 เปรียบเทียบผลการทดลองเมื่อใช้พีเจเออร์แบบต่างๆ ของชุดเอกสาร PDDP J-series	47
4.6 ชุดเอกสาร PDDP K-series ที่ใช้ในการทดลองและผลการทดลอง	53
4.6.1 ชุดเอกสาร PDDP K-series ที่ใช้ในการทดลอง	53
4.6.2 ผลการทดลองของชุดเอกสาร PDDP K-series เมื่อใช้พีเจเออร์ที่ได้ จากการเลือกโดยใช้ต้นไม้วิลิ.....	53
4.6.3 ผลการทดลองของชุดเอกสาร PDDP K-series เมื่อใช้พีเจเออร์แบบ คำเดี่ยว.....	54
4.6.4 ผลการทดลองของชุดเอกสาร PDDP K-series เมื่อใช้พีเจเออร์แบบ คำเดี่ยวและ ไป-แกรม	57
4.6.5 ผลการทดลองของชุดเอกสาร PDDP K-series เมื่อใช้พีเจเออร์แบบ คำเดี่ยวและวลี	58
4.6.6 เปรียบเทียบผลการทดลองเมื่อใช้พีเจเออร์แบบต่างๆ ของชุดเอกสาร PDDP K-series	59
4.7 ชุดเอกสาร Reuters-top10 ที่ใช้ในการทดลองและผลการทดลอง	64
4.7.1 ชุดเอกสาร Reuters-top10 ที่ใช้ในการทดลอง	64

สารบัญ (ต่อ)

	หน้า
4.7.2 ผลการทดลองของชุดเอกสาร Reuters-top10 เมื่อใช้ฟเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลี.....	65
4.7.3 ผลการทดลองของชุดเอกสาร Reuters-top10 เมื่อใช้ฟเจอร์แบบคำเดี่ยว.....	66
4.7.4 ผลการทดลองของชุดเอกสาร Reuters-top10 เมื่อใช้ฟเจอร์แบบคำเดี่ยวและไป-แกรม.....	68
4.7.5 ผลการทดลองของชุดเอกสาร Reuters-top10 เมื่อใช้ฟเจอร์แบบคำเดี่ยวและวลี.....	69
4.7.6 เปรียบเทียบผลการทดลองเมื่อใช้ฟเจอร์แบบต่างๆ ของชุดเอกสาร Reuters-top10.....	71
4.7.7 เปรียบเทียบผลการทดลองเมื่อใช้ฟเจอร์แบบต่างๆ ต่างแบบจำกัดจำนวนของชุดเอกสาร Reuters-top10.....	75
4.8 การเปรียบเทียบผลการทดลองทุกชุดข้อมูล.....	83
 บทที่ 5 สรุปผลการวิจัยและข้อเสนอแนะ.....	 85
 เอกสารอ้างอิง.....	 86
 ภาคผนวก.....	 88
 ภาคผนวก ก. แสดงตัวอย่างชุดเอกสารต่างๆ ที่ใช้ในการทดลอง.....	 89
ก.1 ชุดเอกสาร PDDP.....	89
ก.1.1 ตัวอย่างไฟล์เอกสาร PDDP F-series ในรูปแบบ HTML.....	89
ก.1.2 ตัวอย่างไฟล์เอกสาร PDDP F-series ที่ตัดแท็กออกแล้ว.....	90
ก.1.3 วิธีที่สกัดได้จากตัวอย่างไฟล์เอกสาร PDDP F-series.....	91
ก.2 การสุ่มตัวอย่างเอกสารของชุดเอกสาร PDDP K-series.....	91

สารบัญ (ต่อ)

	หน้า
ก.3 ชุดเอกสาร Reuters-top10	93
ก.3.1 ตัวอย่างไฟล์เอกสาร Reuters-top10ในรูปแบบ XML.....	94
ก.3.2 ตัวอย่างไฟล์เอกสาร Reuters-top10 ที่ตัดแท็กออกแล้ว.....	94
ก.3.3 วิธีที่สกัดได้จากตัวอย่างไฟล์เอกสาร Reuters-top10	94
ก.4 การสุ่มตัวอย่างเอกสารของชุดเอกสาร Reuters-top10	95
ก.5 อัลกอริทึม Phrase-tree-based Feature Reduction	97
ก.6 รายการคำหยุด (Stopwords) ที่ใช้ในการทดลอง.....	98
ภาคผนวก ข. ชุดเอกสาร ฟีเจอร์ของชุดเอกสารต่างๆ ที่ได้จากการทดลอง.....	101
ข.1 ฟีเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลีของชุดเอกสาร F-series	101
ข.2 ฟีเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลีของชุดเอกสาร J-series	123
ข.3 ฟีเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลีของชุดเอกสาร K-series	136
ข.4 ฟีเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลีของชุดเอกสาร Reuters-top10	159
ภาคผนวก ค. ผลงานวิจัยที่เกี่ยวข้องกับวิทยานิพนธ์และได้รับการตีพิมพ์.....	181
ประวัติผู้เขียน.....	187

สารบัญตาราง

ตารางที่	หน้า
2.1 แสดงตัวอย่างเอกสาร เนื้อหาและเทอมที่ได้หลังการเตรียมเอกสาร	12
4.1 แสดงจำนวนเอกสาร F-series แบ่งตามประเภทเอกสาร	32
4.2 แสดงจำนวนเอกสาร F-series ที่ใช้ในการเรียนรู้ และทดสอบในแต่ละชุดการทดลอง	32
4.3 จำนวนพีเจอร์จากต้นไม้ลิ ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง ของชุดเอกสาร PDDP F-series	32
4.4 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เน็ยเรสเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลอง เมื่อใช้พีเจอร์จากต้นไม้ลิของชุดเอกสาร PDDP F-series	33
4.5 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เน็ยเรสเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์จากต้นไม้ลิ ของชุดเอกสาร PDDP F-series	33
4.6 แสดงจำนวนพีเจอร์แบบคำเดี่ยวที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลองของชุดเอกสาร PDDP F-series	34
4.7 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เน็ยเรสเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์แบบคำเดี่ยว ของชุดเอกสาร PDDP F-series	34
4.8 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เน็ยเรสเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์แบบคำเดี่ยว ของชุดเอกสาร F-series	34
4.9 แสดงจำนวนคำเดี่ยวและไบ-แกรมที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง ของชุดเอกสาร PDDP F-series	35
4.10 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เน็ยเรสเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์แบบคำเดี่ยวและไบ-แกรม ของชุดเอกสาร PDDP F-series	35

สารบัญตาราง (ต่อ)

ตารางที่	หน้า
4.11 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการเมื่อใช้ฟิเจอร์แบบคำเดี่ยวและไบ-แกรม ของชุดเอกสาร F-series	36
4.12 แสดงจำนวนฟิเจอร์คำเดี่ยวและวลี ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง ของชุดเอกสาร PDDP F-series	36
4.13 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟิเจอร์คำเดี่ยวและวลี ของชุดเอกสาร PDDP F-series	37
4.14 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการเมื่อใช้ฟิเจอร์คำเดี่ยวและวลี ของชุดเอกสาร F-series	37
4.15 แสดงจำนวนเอกสาร J-series แบ่งตามประเภท.....	42
4.16 แสดงจำนวนเอกสาร J-series ที่ใช้ในการเรียนรู้ และทดสอบในแต่ละชุดการทดลอง	43
4.17 จำนวนฟิเจอร์จากต้นไม้วลี ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง ของชุดเอกสาร PDDP J-series	44
4.18 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลอง เมื่อใช้ฟิเจอร์จากต้นไม้วลีของชุดเอกสาร PDDP J-series	44
4.19 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟิเจอร์จากต้นไม้วลี ของชุดเอกสาร PDDP J-series	44
4.20 แสดงจำนวนฟิเจอร์แบบคำเดี่ยวที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง ของชุดเอกสาร PDDP J-series	45
4.21 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟิเจอร์แบบคำเดี่ยว ของชุดเอกสาร PDDP J-series	45

สารบัญตาราง (ต่อ)

ตารางที่	หน้า
4.22 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์แบบคำเดี่ยว ของชุดเอกสาร J-series	45
4.23 แสดงจำนวนคำเดี่ยวและไป-แกรมที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง ของชุดเอกสาร PDDP J-series	46
4.24 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์แบบคำเดี่ยวและไป-แกรม ของชุดเอกสาร PDDP J-series	46
4.25 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการเมื่อใช้พีเจอร์แบบคำเดี่ยวและไป-แกรม ของชุดเอกสาร J-series	47
4.26 แสดงจำนวนพีเจอร์คำเดี่ยวและวลี ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง ของชุดเอกสาร PDDP J-series	47
4.27 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์คำเดี่ยวและวลี ของชุดเอกสาร PDDP J-series	48
4.28 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการเมื่อใช้พีเจอร์คำเดี่ยวและวลี ของชุดเอกสาร J-series.....	48
4.29 แสดงจำนวนเอกสาร K-series ที่ใช้ในการเรียนรู้ และทดสอบในแต่ละชุดการทดลอง	54
4.30 จำนวนพีเจอร์จากต้นไม้วลี ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง ของชุดเอกสาร PDDP K-series	55
4.31 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลอง เมื่อใช้พีเจอร์จากต้นไม้วลีของชุดเอกสาร PDDP K-series	55

สารบัญตาราง (ต่อ)

ตารางที่	หน้า
4.32 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟัเจอร์จากต้นไม้วลี ของชุดเอกสาร PDDP K-series	55
4.33 แสดงจำนวนฟัเจอร์แบบคำเดี่ยวที่ได้จากชุดเอกสารเรียนรู ในแต่ละชุดการทดลองของชุดเอกสาร PDDP K-series	56
4.34 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟัเจอร์แบบคำเดี่ยว ของชุดเอกสาร PDDP K-series	56
4.35 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟัเจอร์แบบคำเดี่ยว ของชุดเอกสาร K-series	56
4.36 แสดงจำนวนคำเดี่ยวและไบ-แกรมที่ได้จากชุดเอกสารเรียนรู ในแต่ละชุดการทดลอง ของชุดเอกสาร PDDP K-series	57
4.37 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟัเจอร์แบบคำเดี่ยวและไบ-แกรม ของชุดเอกสาร PDDP K-series	57
4.38 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการเมื่อใช้ฟัเจอร์แบบคำเดี่ยวและไบ-แกรม ของชุดเอกสาร K-series	58
4.39 แสดงจำนวนฟัเจอร์คำเดี่ยวและวลี ที่ได้จากชุดเอกสารเรียนรู ในแต่ละชุดการทดลอง ของชุดเอกสาร PDDP K-series	58
4.40 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟัเจอร์คำเดี่ยวและวลี ของชุดเอกสาร PDDP K-series	59
4.41 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิชีเค-เนียร์เสนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการเมื่อใช้ฟัเจอร์คำเดี่ยวและวลี ของชุดเอกสาร K-series.....	59

สารบัญตาราง (ต่อ)

ตารางที่	หน้า
4.42 แสดงจำนวนเอกสาร Reuters-top10 ที่ใช้ในการเรียนรู้ และทดสอบในแต่ละชุดการทดลอง	65
4.43 จำนวนฟีเจอร์จากต้นไม้วลี ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง ของชุดเอกสาร Reuters-top10	66
4.44 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลอง เมื่อใช้ฟีเจอร์จากต้นไม้วลีของชุดเอกสาร Reuters-top10	66
4.45 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์จากต้นไม้วลี ของชุดเอกสาร Reuters-top10	67
4.46 แสดงจำนวนฟีเจอร์แบบคำเดี่ยวที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง ของชุดเอกสาร Reuters-top10	67
4.47 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์แบบคำเดี่ยว ของชุดเอกสาร Reuters-top10.....	67
4.48 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์แบบคำเดี่ยว ของชุดเอกสาร K-series	68
4.49 แสดงจำนวนคำเดี่ยวและไบ-แกรมที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง ของชุดเอกสาร Reuters-top10	68
4.50 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์แบบคำเดี่ยวและไบ-แกรม ของชุดเอกสาร Reuters-top10	69
4.51 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการเมื่อใช้ฟีเจอร์แบบคำเดี่ยวและไบ-แกรม ของชุดเอกสาร K-series	69

สารบัญตาราง (ต่อ)

ตารางที่	หน้า
4.52 แสดงจำนวนพีเจอร်ค่าเดียวและวลี ที่ได้จากชุดเอกสารเรียนรู้ในแต่ละชุดการทดลอง ของชุดเอกสาร Reuters-top10	70
4.53 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร်ค่าเดียวและวลี ของชุดเอกสาร Reuters-top10.....	70
4.54 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร်ค่าเดียวและวลี ของชุดเอกสาร Reuters-top10.....	70
ก.1 แสดงจำนวนเอกสารตัวอย่าง K-series แบ่งตามประเภทเอกสาร	91
ก.2 แสดงจำนวนเอกสารของชุดเอกสาร Reuters-top10 แบ่งตามประเภทเอกสาร.....	94

สารบัญรูป

รูปที่	หน้า
2.1 การจำแนกประเภทเอกสาร.....	6
2.2 กระบวนการในการจำแนกประเภทเอกสาร	6
2.3 การเลือกฟีเจอร์ในระดับ Local และ Global	9
2.4 การเลือกเซตย่อยของฟีเจอร์ที่ดีที่สุด	9
2.5 อัลกอริทึมการจำแนกประเภทเอกสารแบบ kNN.....	14
3.1 กระบวนการของการจำแนกประเภทเอกสาร โดยทำการลดฟีเจอร์โดยใช้ต้นไม้ตัดสินใจ ..	19
3.2 ขั้นตอนการลดฟีเจอร์โดยใช้ต้นไม้ตัดสินใจ.....	21
3.3 ตัวอย่างต้นไม้ตัดสินใจที่ได้จากวลี “ <i>k-nearest neighbor</i> ” และ “ <i>effici categor algorithm</i> ”	22
3.4 ตัวอย่างต้นไม้ตัดสินใจจากรูปที่ 3.3 ที่ให้คะแนนแล้ว	23
3.5 ตัวอย่างต้นไม้ตัดสินใจจากรูป 3.4 ที่ผ่านการเลือกฟีเจอร์แล้ว.....	23
3.6 ตัวอย่างต้นไม้ตัดสินใจของคลาส c_2 ที่ผ่านการเลือกฟีเจอร์ระดับ Localแล้ว.....	24
3.7 ต้นไม้ตัดสินใจจากการรวมกันของต้นไม้ในรูปที่ 3.5 และ 3.6	24
3.8 ต้นไม้ตัดสินใจที่ได้หลังจากทำการลดฟีเจอร์ จากต้นไม้ตัดสินใจในรูปที่ 3.7	25
3.9 อัลกอริทึมการลดฟีเจอร์โดยใช้ต้นไม้ตัดสินใจ.....	26
4.1 แสดงส่วนที่ใช้วัดค่า Precision และ Recall ในการร้องขอสารสนเทศ	28
4.2 การทำ 4-fold cross validation ในชุดการทดลองลำดับที่ 1 (fold 1)	30
4.3 แผนภูมิแท่งเปรียบเทียบจำนวนฟีเจอร์ที่ใช้ในแต่ละแบบของชุดเอกสาร F-series ...	37
4.4 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 1 .	38
4.5 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 2 .	38
4.6 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 3 .	38
4.7 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 4 .	39
4.8 กราฟเปรียบเทียบอัตราความถูกต้องเฉลี่ยของทุกชุดการทดลองในแต่ละแบบของชุดเอกสาร F-series	39
4.9 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุดในแต่ละแบบของชุดเอกสาร F-series	39
4.10 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 1	40
4.11 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 2	40

สารบัญรูป (ต่อ)

รูปที่	หน้า
4.12 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 3	40
4.13 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 4	41
4.14 กราฟเปรียบเทียบ ค่า F-measure เฉลี่ยในแต่ละแบบของชุดเอกสาร F-series	41
4.15 แผนภูมิแท่งเปรียบเทียบ ค่า F-measure สูงสุดในแต่ละแบบของชุดเอกสาร F-series	41
4.16 แผนภูมิแท่งเปรียบเทียบจำนวนพีเจอร์ที่ใช้ในแต่ละแบบของชุดเอกสาร J-series ...	48
4.17 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 1 .	49
4.18 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 2 .	49
4.19 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 3 .	49
4.20 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 4 .	50
4.21 กราฟเปรียบเทียบอัตราความถูกต้องเฉลี่ยของทุกชุดการทดลองในแต่ละแบบของชุดเอกสาร J-series	50
4.22 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุดในแต่ละแบบของชุดเอกสาร J-series	50
4.23 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 1	51
4.24 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 2	51
4.25 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 3	51
4.26 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 4	52
4.27 กราฟเปรียบเทียบ ค่า F-measure เฉลี่ยในแต่ละแบบของชุดเอกสาร J-series	52
4.28 แผนภูมิแท่งเปรียบเทียบ ค่า F-measure สูงสุดในแต่ละแบบของชุดเอกสาร J-series	52
4.29 แผนภูมิแท่งเปรียบเทียบจำนวนพีเจอร์ที่ใช้ในแต่ละแบบของชุดเอกสาร K-series .	60
4.30 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 1.	60
4.31 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 2.	60
4.32 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 3.	61
4.33 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 4.	61
4.34 กราฟเปรียบเทียบอัตราความถูกต้องเฉลี่ยของทุกชุดการทดลองในแต่ละแบบของชุดเอกสาร K-series	61

สารบัญรูป (ต่อ)

รูปที่	หน้า
4.35 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุดในแต่ละแบบของชุดเอกสาร K-series	62
4.36 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 1	62
4.37 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 2	62
4.38 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 3	63
4.39 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 4	63
4.40 กราฟเปรียบเทียบ ค่า F-measure เฉลี่ยในแต่ละแบบของชุดเอกสาร K-series	63
4.41 แผนภูมิแท่งเปรียบเทียบ ค่า F-measure สูงสุดในแต่ละแบบของชุดเอกสาร K-series	64
4.42 แผนภูมิแท่งเปรียบเทียบจำนวนพีเจอร์ที่ใช้ในแต่ละแบบของชุดเอกสาร Reuters-top10	71
4.43 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 1	71
4.44 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 2	71
4.45 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 3	72
4.46 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 4	72
4.47 กราฟเปรียบเทียบอัตราความถูกต้องเฉลี่ยของทุกชุดการทดลองในแต่ละแบบของชุดเอกสาร Reuters-top10	72
4.48 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุดในแต่ละแบบของชุดเอกสาร Reuters-top10	73
4.49 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 1	73
4.50 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 2	73

สารบัญญรูป (ต่อ)

รูปที่	หน้า
4.51 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 3	74
4.52 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 4	74
4.53 กราฟเปรียบเทียบ ค่า F-measure เฉลี่ยในแต่ละแบบของชุดเอกสาร Reuters-top10 ..	74
4.54 แผนภูมิแท่งเปรียบเทียบ ค่า F-measure สูงสุด ในแต่ละแบบของชุดเอกสาร Reuters-top10	75
4.55 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 1250 ตัว ของชุดเอกสาร Reuters-top10	76
4.56 แผนภูมิแท่งเปรียบเทียบค่า F-measure สูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 1250 ตัว ของชุดเอกสาร Reuters-top10	76
4.57 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 1500 ตัว ของชุดเอกสาร Reuters-top10	77
4.58 แผนภูมิแท่งเปรียบเทียบค่า F-measure สูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 1500 ตัว ของชุดเอกสาร Reuters-top10	77
4.59 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 1750 ตัว ของชุดเอกสาร Reuters-top10	78
4.60 แผนภูมิแท่งเปรียบเทียบค่า F-measure สูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 1750 ตัว ของชุดเอกสาร Reuters-top10	78
4.61 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 2000 ตัว ของชุดเอกสาร Reuters-top10	79
4.62 แผนภูมิแท่งเปรียบเทียบค่า F-measure สูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 2000 ตัว ของชุดเอกสาร Reuters-top10	79
4.63 กราฟเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้ฟิเจอร์ในแต่ละแบบ โดยมีการ จำกัดจำนวนของฟิเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 1	80
4.64 กราฟเปรียบเทียบอัตราค่า F-measure สูงสุด เมื่อใช้ฟิเจอร์ในแต่ละแบบ โดยมีการ จำกัดจำนวนของฟิเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 1	80

สารบัญรูป (ต่อ)

รูปที่	หน้า
4.65 กราฟเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 2	80
4.66 กราฟเปรียบเทียบอัตราค่า F-measure สูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 2	81
4.67 กราฟเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 3	81
4.68 กราฟเปรียบเทียบอัตราค่า F-measure สูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 3	81
4.69 กราฟเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 4	82
4.70 กราฟเปรียบเทียบอัตราค่า F-measure สูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 4	82
4.71 แผนภูมิแท่งเปรียบเทียบค่าอัตราความถูกต้องเฉลี่ย เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10.....	82
4.72 แผนภูมิแท่งเปรียบเทียบอัตราค่า F-measure เฉลี่ย เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10.....	83
4.73 แผนภูมิแท่งเปรียบเทียบจำนวนพีเจอร์แต่ละแบบของทุกชุดข้อมูล	83
4.74 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องเมื่อใช้พีเจอร์แต่ละแบบของทุกชุดข้อมูล.....	84
4.75 แผนภูมิแท่งเปรียบเทียบค่า F-measure เมื่อใช้พีเจอร์แต่ละแบบของทุกชุดข้อมูล	84
ก.1 ตัวอย่างเว็บเพจของเอกสารตัวอย่าง PDDP F-series.....	88
ก.2 แสดงชั้นภูมิค่าและจำนวนเอกสารในแต่ละประเภทเอกสารของชุดเอกสาร K-series	90
ก.3 แสดงสัดส่วนระหว่างจำนวนเอกสารทั้งหมดกับเอกสารตัวอย่างในแต่ละประเภทของชุดเอกสาร K-series	90
ก.4 แสดงชั้นภูมิค่าและจำนวนเอกสารในแต่ละประเภทเอกสารของชุดเอกสาร Reuters-top10	93

สารบัญรูป (ต่อ)

รูปที่	หน้า
ก.5 แสดงสัดส่วนระหว่างจำนวนเอกสารทั้งหมดกับเอกสารตัวอย่างในแต่ละประเภท ของชุดเอกสาร Reuters-top10.....	94

บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

ในปัจจุบันมีการใช้คอมพิวเตอร์ และอินเทอร์เน็ตกันอย่างแพร่หลาย ทำให้ข้อมูลที่อยู่ในรูปแบบอิเล็กทรอนิกส์มีปริมาณเพิ่มขึ้นอย่างรวดเร็ว โดยเฉพาะข้อมูลที่อยู่ในรูปแบบเอกสารหรือข้อความ เช่น ข้อมูลงานวิจัย ข้อมูลข่าว และหนังสืออิเล็กทรอนิกส์ ดังนั้นจึงมีความต้องการในการจัดการกับข้อมูลเหล่านี้ เพื่อให้สามารถนำมาใช้ประโยชน์ได้อย่างมีประสิทธิภาพ

การจำแนกประเภทเอกสาร (Text Categorization หรือ Text Classification) เป็นการจัดเอกสารตามลักษณะเนื้อหาให้สอดคล้องกับประเภทเอกสารต่างๆ ที่มีอยู่ ทำให้เอกสารสามารถถูกจัดเก็บอย่างมีระบบ สืบค้นได้ง่าย ซึ่งการจำแนกประเภทเอกสารนั้นยังได้ถูกนำมาประยุกต์ใช้ในงานต่างๆ เช่น การกรองเอกสาร กลไกการค้นหาหัวข้อ การจำแนกประเภทเว็บ การสร้างดัชนีแบบอัตโนมัติในระบบการเทคนิคค้นคืนสารสนเทศ

ส่วนที่สำคัญอย่างหนึ่งในการจำแนกประเภทเอกสาร คือ การแทนเอกสาร (Document Representation) ซึ่งทำให้เอกสารอยู่ในรูปแบบที่จะถูกจำแนกประเภทต่อไป การแทนเอกสารที่ใช้กันอย่างแพร่หลายคือ การแทนเอกสารโดยใช้คำเดี่ยว (Single words) เป็นฟีเจอร์ (Feature) เรียกว่า Bags of word [1] ซึ่งแต่ละเอกสารจะถูกแทนเป็นเวกเตอร์ โดยแต่ละไอเทมในเวกเตอร์คือค่าที่เป็นตัวเลขสอดคล้องกับคำนั้นๆ ซึ่งอาจจะเป็นความถี่ของคำในเอกสารนั้น หรือเป็นค่าไบนารี หรือเป็นค่าน้ำหนักต่างๆ โดยเวกเตอร์ของเอกสารทั้งหมด เรียกว่า Vector Space Model [2] ซึ่งมีขนาดหรือมิติคือ จำนวนเอกสาร \times จำนวนฟีเจอร์

การใช้คำเดี่ยวนี้มีข้อดีหลายอย่างด้วยกัน ยกตัวอย่างเช่น ทำให้ขาดความสัมพันธ์ระหว่างคำ เช่น ลำดับของเทอม และคำๆ หนึ่งยังสามารถมีได้หลายความหมายถ้าปรากฏในเนื้อหาที่ต่างกัน วิธีหนึ่งในการแก้ปัญหาดังกล่าว คือ การใช้วลี (Phrases) เป็นฟีเจอร์แทนการใช้คำเดี่ยว วลีสามารถลดปัญหาความไม่แน่นอนของความหมายของคำเดี่ยวได้ ยกตัวอย่างเช่น คำว่า “Java” มี 2 ความหมายด้วยกัน คือ ภาษาโปรแกรมชนิดหนึ่ง และเป็นชื่อเกาะของประเทศอินโดนีเซีย แต่ถ้าเราใช้วลีเป็น “Java Script” หรือ “Java Island” ก็จะได้ความหมายที่ชัดเจนมากกว่า

ดังนั้น จึงมีการนำวลีมาใช้เป็นฟีเจอร์ในการแทนเอกสาร แต่ในการศึกษาการจำแนกประเภทเอกสาร โดยใช้วลีเป็นฟีเจอร์เพียงอย่างเดียว นั้น ผลลัพธ์ที่ได้กลับมีความถูกต้องลดลงน้อยกว่าการใช้คำเดี่ยว ซึ่งเหตุผลหลักก็คือ มีจำนวนวลีต่างๆ มาก แต่ว่าแต่ละวลีกลับมีความถี่ที่น้อยมาก และวลีดังกล่าวยังมีความซ้ำซ้อนกันอยู่มาก คือมีวลีจำนวนมากที่มีความหมายเหมือนกัน จึง

มีผู้นำคำเดี่ยวและวลี มาใช้เป็นฟีเจอร์ร่วมกัน ซึ่งช่วยเพิ่มความถูกต้องในการจำแนกประเภทเอกสารมากขึ้น

วลีที่ใช้ในเป็นฟีเจอร์นั้นแบ่งได้ 2 แบบ คือ วลีทางภาษา (Syntactic Phrase) วลีในลักษณะนี้ จะถูกต้องตามหลักไวยากรณ์ของภาษา และใช้วิธีการทางภาษาศาสตร์ (NLP) ในการแยกวลีออกมา ส่วนอีกแบบหนึ่งคือ วลีทางสถิติ (Statistical Phrase) หมายถึง เซต หรือ ลำดับของคำใดๆ ที่ปรากฏต่อเนื่องกันในเอกสาร โดยมีค่านัยสำคัญทางสถิติ ตัวอย่างวลีประเภทนี้เช่น การสร้าง n -gram ซึ่งคือ ลำดับของคำใด ๆ ที่ปรากฏต่อเนื่องกันในเอกสาร

ปัญหาของวลีทางภาษา คือ ในการสกัดวลี (Phrase Extraction) ทางภาษานั้นทำได้ยากกว่าวลีทางสถิติ และวลีทางภาษาที่ถูกตามหลักไวยากรณ์นั้นจะมีความหมายแคบเกินไป เช่น วลี “*k-nearest neighbor categorization algorithm*” จะปรากฏในแต่ละเอกสารน้อยมาก แต่ถ้าเราใช้วลีย่อยเป็น “*k-nearest neighbor*” หรือ “*categorization algorithm*” จะมีการกระจายของคำในการแทนเอกสารดีกว่า เพราะมีโอกาสในการปรากฏมากกว่าในแต่ละเอกสารทำให้คำที่ถูกแทนในการแทนเอกสารมีค่าน้อย แต่ถ้าวลีนั้นสั้นไป ก็จะทำให้เกิดปัญหาคัดเลือกกับการใช้คำเดี่ยว คือ มีความสามารถในการแบ่งแยกประเภทต่ำ โดยเฉพาะอย่างยิ่งเมื่อประเภทของข้อมูลมีความใกล้เคียงกันมาก เช่น วลี “*drama player voting*” อยู่ในข่าวบันเทิง ส่วนวลี “*football player voting*” อยู่ในข่าวกีฬา ถ้าเราใช้ฟีเจอร์เป็น “*player voting*” ก็ทำให้ประสิทธิภาพในการจำแนกว่าอยู่ในข่าวประเภทไหนลดลง

การใช้วลีทางสถิติเป็นการแก้ปัญหาอย่างหนึ่งของการใช้วลีทางภาษา เพราะไม่ต้องคำนึงถึงหลักไวยากรณ์ทางภาษา เช่น การสร้าง n -gram ซึ่งคือ ลำดับของคำใดๆ ที่ปรากฏต่อเนื่องกันในเอกสาร โดยส่วนใหญ่จะใช้ $n=2$ เรียกว่า Bi-gram แต่ปัญหาก็คือ วลีที่สร้างได้มีจำนวนมาก ที่เป็นวลีขยะ คือ วลีที่มีอำนาจในการจำแนกประเภทต่ำ และในการสร้าง n -gram นั้น จำเป็นที่จะต้องกำหนดขนาดสูงสุดของวลีไว้ คือ n ทำให้บางครั้งมีวลีที่มีความสามารถในการจำแนกประเภทที่ดี แต่มีขนาดมากกว่า n ก็จะไม่ถูกใช้

วิธีการแก้ปัญหาที่กล่าวมาก็คือ จะต้องมียุทธวิธีที่ดีในการสร้างวลี และต้องทำการเลือกเอาเฉพาะวลีที่สามารถแบ่งแยกประเภทของเอกสารได้ดีเท่านั้น โดยทำการคัดเลือกฟีเจอร์ (Feature Selection หรือ Feature Reduction) [3] เพื่อช่วยลดขนาดของฟีเจอร์ที่ใช้ และช่วยเพิ่มประสิทธิภาพการจำแนกประเภทเอกสาร

1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา

วิทยานิพนธ์ฉบับนี้ มีวัตถุประสงค์เพื่อพัฒนาวิธีการ และอัลกอริทึมในการสร้างฟีเจอร์วลี และเลือกวลี เพื่อลดขนาดของฟีเจอร์วลีที่จะถูกนำมาใช้แทนเอกสาร สำหรับการจำแนกประเภทเอกสาร โดยประยุกต์ใช้ทฤษฎีการสร้างฟีเจอร์ และการคัดเลือกฟีเจอร์ ร่วมกับ โครงสร้างข้อมูล

แบบต้นไม้ ที่ช่วยทำให้การจำแนกประเภทเอกสารมีประสิทธิภาพดียิ่งขึ้น ทั้งในแง่ของขนาดโดยการลดมิติเพื่อการทำงานที่เร็วขึ้น และในแง่ของความถูกต้องในการจำแนกประเภทเอกสาร

1.3 ทฤษฎีหรือแนวคิดที่ใช้ในการวิจัย

งานวิจัยนี้ได้นำทฤษฎีและเทคนิคต่างๆมาประยุกต์ใช้ประกอบด้วย

1. เทคนิคการจำแนกประเภทเอกสาร ตั้งแต่การเตรียมเอกสาร การสร้างแบบจำลองในการจำแนกประเภทไปจนถึงการวัดประสิทธิภาพ
2. เทคนิคการแทนเอกสาร
3. เทคนิคการเลือกฟีเจอร์
4. เทคนิคการดำเนินการ โครงสร้างข้อมูลแบบต้นไม้

1.4 ขอบเขตการวิจัย

ในการวิจัยนี้จะออกแบบและพัฒนาวิธีการและอัลกอริทึมที่ใช้โครงสร้างต้นไม้ ร่วมกับเทคนิคการเลือกฟีเจอร์ เพื่อคัดเลือกฟีเจอร์ที่เป็นวลี ที่มีประสิทธิภาพ และลดจำนวนของฟีเจอร์ที่เป็นวลีที่ใช้ในการแทนเอกสาร

1.5 ขั้นตอนของการศึกษา

ในขั้นตอนของการศึกษานี้ ได้แสดงลำดับการทำงานตั้งแต่เริ่มต้นจนถึงสิ้นสุดการทำงาน ดังรายละเอียดต่อไปนี้

- 1.5.1 ศึกษาทฤษฎีและงานวิจัยจากเอกสารบทความต่างๆ ที่เกี่ยวข้องกับการทำงานวิจัย
- 1.5.2 กำหนด หัวข้อ วัตถุประสงค์ และขอบเขตการทำงานวิจัย
- 1.5.3 ออกแบบและวิเคราะห์อัลกอริทึมใหม่
- 1.5.4 พัฒนาโปรแกรม โดยใช้ซอฟต์แวร์ MATLAB รุ่น 6.5 พร้อมทั้งแก้ไขข้อผิดพลาด และทดสอบการทำงานของอัลกอริทึม
- 1.5.5 เตรียมข้อมูลที่ใช้งานจริง เพื่อนำมาทดสอบการทำงานของอัลกอริทึม
- 1.5.6 ทดลองกับข้อมูลที่ใช้งานจริง พร้อมทั้งวัดประสิทธิภาพการทำงานของอัลกอริทึม
- 1.5.7 รวบรวมผลการทดลองจากการทำงานของอัลกอริทึม
- 1.5.8 วิเคราะห์และสรุปผลการทดลอง
- 1.5.9 เรียบเรียงเอกสารประกอบวิทยานิพนธ์

งานวิจัยนี้มุ่งเน้นศึกษาในส่วนของฟิวเจอร์ที่เป็นวัสดุที่ใช้ในการแทนเอกสาร โดยพยายามสร้าง และเลือกใช้ฟิวเจอร์ที่เป็นวัสดุ เฉพาะที่มีประสิทธิภาพเพื่อลดขนาดของฟิวเจอร์วัสดุที่ใช้ และนำโครงสร้างต้น ไม้มาประยุกต์ใช้ในการทำงาน

ในบทที่ 2 จะกล่าวถึง ทฤษฎีต่างๆ ที่เกี่ยวข้องกับการทำงานวิจัยที่นำเสนอนี้

บทที่ 2

ทฤษฎีพื้นฐาน และงานวิจัยที่เกี่ยวข้อง

ในหัวข้อนี้จะกล่าวถึงทฤษฎีพื้นฐานต่างๆ ที่เกี่ยวข้องในการวิจัย พื้นฐานของ การจำแนกประเภทเอกสาร การลดฟีดเจอร์ และแสดงถึงการจำแนกประเภทเอกสาร โดยใช้ฟีดเจอร์ในลักษณะต่างๆ พร้อมทั้งแสดงถึงงานวิจัยที่เกี่ยวข้อง และปัญหาที่เกิดขึ้น

2.1 การจำแนกประเภทเอกสาร (Text Categorization)

การจำแนกประเภทเอกสาร คือ การจัดเอกสารไปยังประเภทต่างๆ ตามลักษณะของเนื้อหา นั้น ๆ ว่าสอดคล้องกับเอกสารในประเภทไหนที่มีอยู่ โดยการจำแนกประเภทเป็นงานหนึ่งของดาต้า ไมนิ่ง (Data Mining) ที่ทำการค้นหาสารสนเทศที่มีอยู่จากข้อมูลที่มีปริมาณมาก ซึ่งการจำแนกประเภทอยู่ในส่วนของการทำนาย (Prediction) โดยหมายถึง การนำตัวอย่างจากประสบการณ์ที่ผ่านมาที่เรา รู้คำตอบ แน่ชัดแล้ว มาใช้ในกรณีในอนาคต [24] ในการจำแนกประเภทตัวอย่างที่เรา รู้แน่ชัดแล้วก็คือ เอกสารที่ถูกให้ป้ายชื่อของประเภทไว้ (Label) การจำแนกประเภทแตกต่างจากการจัดกลุ่ม (Clustering) ตรงที่การจัดกลุ่ม เราไม่รู้เลยว่าข้อมูลอยู่ในประเภทไหน และแบ่งได้กี่ประเภท

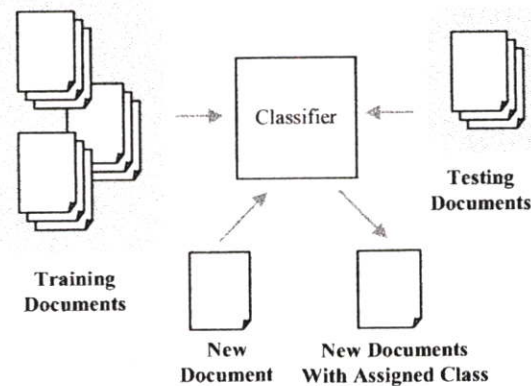
Sebastiani [4] ได้นิยามการจำแนกประเภทเอกสารไว้ว่า การจำแนกประเภทเอกสาร คือ งานของการให้ค่าบวกลบกับแต่ละคู่ระหว่างเอกสารกับคลาส $(d_j, c_i) \in D \times C$ โดย D คือ เซตของเอกสาร และ C คือ เซตของคลาส

การจำแนกประเภทถูกนำมาประยุกต์ใช้ในหลายๆ งานด้วยกัน [4] เช่น

1. Document Organization: เช่น ในงานหนังสือพิมพ์ มีความจำเป็นที่จะต้องการจำแนกประเภทของข่าวที่ผ่านมา เพื่อความง่ายในการค้นหาในอนาคต ในกรณีที่เหตุการณ์ที่เกิดขึ้นใหม่มีเกี่ยวข้องกับข่าวที่ผ่านมา ประเภทที่เป็นไปได้ เช่น “Sport”, “Entertainment”
2. การกรองเอกสาร (Text Filtering): เช่น ในการกรอง e-mail โดยแบ่ง e-mail เป็น 2 ประเภท คือ e-mail ที่เราสนใจ กับ e-mail ที่เราไม่สนใจ หรือ e-mail ที่เป็นขยะ กับ e-mail ที่ไม่เป็นขยะ
3. การแบ่งประเภทแบบมีลำดับชั้นของเว็บเพจ

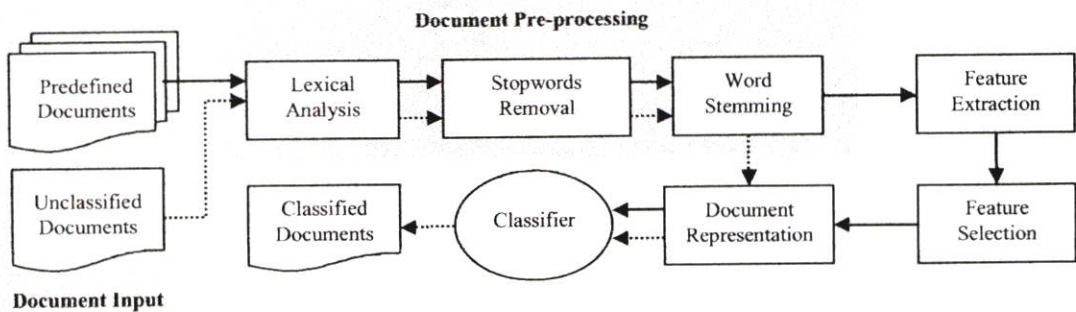
การจำแนกประเภทเอกสาร ประกอบด้วยกระบวนการหลักๆ ดังนี้ คือ การเตรียมเอกสาร การแทนเอกสาร การสร้างตัวจำแนกประเภท (Classifier) ซึ่งในขั้นตอนนี้ ชุดเอกสารที่อยู่ในประเภทต่างๆ จะถูกแบ่งเป็น 2 ส่วน เพื่อใช้ในการเรียนรู้ในการสร้างตัวจำแนกประเภท และอีกชุด

เพื่อใช้ในการทดสอบตัวจำแนกประเภท โดยการวัดประสิทธิภาพในการจำแนกประเภทเอกสารของตัวจำแนกประเภท และกระบวนการสุดท้ายก็คือ นำเอกสารใหม่มาจำแนกประเภทโดยใช้ตัวจำแนกประเภทที่สร้างไว้ ดังแสดงในรูปที่ 2.1



รูปที่ 2.1 การจำแนกประเภทเอกสาร

กระบวนการในการจำแนกประเภทเอกสารแสดงดังรูปที่ 2.2 เส้นทึบแทนการไหลของข้อมูลในขั้นตอนการสร้างตัวจำแนกประเภท โดยเอกสารที่เรารู้คลาสแล้วจะผ่านกระบวนการเตรียมเอกสาร จากนั้นผ่านกระบวนการการสกัดและเลือกฟีเจอร์ และการแทนเอกสาร เพื่อสร้างเป็นตัวจำแนกประเภท ส่วนเส้นประหมายถึงเส้นทางการไหลของข้อมูลในขั้นตอนการจำแนกประเภท โดยเมื่อมีเอกสารที่เรายังไม่ทราบคลาสเข้ามา ก็ผ่านกระบวนการเตรียมเอกสาร และการแทนเอกสาร เพื่อนำเข้าสู่ตัวจำแนกประเภท ทำให้ทราบประเภทของเอกสารนั้น



รูปที่ 2.2 กระบวนการในการจำแนกประเภทเอกสาร

2.2 การเตรียมเอกสาร (Document Preprocessing)

ขั้นตอนแรกของการจำแนกประเภทเอกสาร (Text Categorization) คือการเตรียมเอกสารให้อยู่ในรูปแบบที่จะนำเข้าสู่ขั้นตอนการจำแนกประเภท ได้ ซึ่งจะมีการดำเนินการเปลี่ยนรูปแบบเอกสาร (Text Transformations หรือ Text Operations) [6] 3 ลักษณะ ดังต่อไปนี้

1. การวิเคราะห์คำ (Lexical Analysis) หรือ โทเคนไนเซชัน (Tokenization): เอกสารจะถูกพิจารณาเป็นสายอักขระ (String) อักขระที่ไม่ใช่ตัวอักษร เช่น ช่องว่าง (Space) บรรทัดใหม่ (New Line) แท็บ (Tab) และเครื่องหมายวรรคตอนต่างๆ จะถูกเอาออก โดยเก็บส่วนที่เหลือไว้ และเปลี่ยนให้เป็นตัวพิมพ์เล็ก (Lower Case) ให้หมด รวมทั้งการเอาแท็ก (Tags) ต่างๆ ของพวกเอกสารที่เป็นเว็บออกด้วย

2. การกำจัด Stopwords (Elimination of Stopwords) คือ การตัดคำที่มีอำนาจในการแบ่งแยกได้น้อยออกไป หรือ Function Word เช่น คำนำหน้า (Articles) เช่น a, an, the คำบุพบท (Prepositions) เช่น in, on, of คำสันธาน (Conjunctions) เช่น and, or, but, if คำสรรพนาม (Pronouns) เช่น I, you, them, it คำต่างๆ ที่เป็นไปได้ เช่น make, thing, similar, all, already, always, nobody โดยรายการคำหยุด สามารถดาวน์โหลดได้ที่ [22]

3. การทำให้อยู่ในรูปแบบของรากคำ (Word Stemming) คือ การแทนที่รูปแปร (Variant) ต่างๆ ของคำ ด้วยรากของคำเพียงคำเดียว รูปแปรประกอบด้วย พหูพจน์ รูปแบบของคำกริยาที่เติมท้ายด้วย “ing” หน่วยคำเติมหลังบุคคลที่สาม หน่วยคำเติมหลังอดีตกาล (Past Tense) เป็นต้น ตัวอย่างเช่น คำว่า connect เป็นรากของคำต่อไปนี้ connects, connected, connecting, connection เป็นต้น ซึ่งคำต่างๆ เหล่านี้มีความหมายที่คล้ายคลึงกัน โดยมาตรฐานที่ใช้ส่วนใหญ่ คือ อัลกอริทึมการสเต็มมิงของพอร์เตอร์ (Porter’s Stemming Algorithm) [7]

2.3 การสกัดฟีเจอร์ (Feature Extraction)

การสกัดฟีเจอร์ คือ การเลือกว่าเราจะใช้ทอมอะไรที่เหมาะสมเป็นฟีเจอร์ เพื่อใช้แทนเอกสาร หรือเรียกว่าการทำดัชนี (Indexing) โดยฟีเจอร์ ที่เราใช้ อาจจะเป็นคำเดี่ยว หรือ วลี โดยในการทำดัชนีนั้นมี 2 วิธีด้วยกัน [4] คือ

1. Derived Indexing: ฟีเจอร์จากวิธีนี้จะได้มาจากเนื้อหาของเอกสาร โดยตรง โดยฟีเจอร์ อาจจะเป็น คำ วลี หรือ n-grams เป็นต้น
2. Assigned Indexing: ฟีเจอร์จากวิธีนี้จะได้มาจากข้อมูลภายนอก เช่น รายการของคำศัพท์ หรือมีการขยายคำที่ใช้เป็นฟีเจอร์ด้วยคำพ้อง (Synonym) จาก WordNet [23] เป็นต้น

โดยฟีเจอร์ที่ใช้กันโดยทั่วไป คือ คำ และวลี โดยมีรายละเอียดดังนี้

2.3.1 คำเดี่ยว (Single Word)

การใช้คำเป็นฟีเจอร์นี้นิยมใช้กันโดยทั่วไป โดยจะเรียกเซตของฟีเจอร์ที่เป็นคำนี้ว่า Bag of Word (BOW) [1] ซึ่งเป็นการใช้คำสำคัญ (Single Word) ที่ได้จากแต่ละเอกสารมาเป็นฟีเจอร์ ใน

เอกสารภาษาอังกฤษซึ่งเป็นภาษาที่เราใช้ในงานวิจัยนี้ สามารถสกัดคำออกมาได้ง่าย เพราะคำในภาษาอังกฤษจะถูกแบ่งแยกด้วยช่องว่างหรือเครื่องหมายต่างๆ ซึ่งต่างจากภาษาไทยที่เราไม่สามารถรู้ขอบเขตของคำอย่างชัดเจนต้องอาศัยพจนานุกรมมาช่วย

2.3.2 วลี (Phrase)

Bag of Phrase (BOP) ใช้วลี (Phrase) ซึ่งคือลำดับของคำ หรือกลุ่มของคำ (Multi-word) เป็นฟีเจอร์ โดยวลีที่ใช้ในการแทนเอกสาร แบ่งได้ 2 ลักษณะ คือ

1. วลีทางภาษา (Syntactic Phrase) วลีในลักษณะนี้ จะถูกต้องตามหลักไวยากรณ์ของภาษา และถูกกำหนดโดยใช้วิธีการทางภาษาศาสตร์ (NLP) ต่างๆ โดย วลีในที่นี้หมายถึง หน่วยของเนื้อหาที่โดยทั่วไปจะใหญ่กว่าคำ แต่เล็กกว่าประโยคเต็ม ตัวอย่างของวลีที่เป็นคำนาม เช่น “nuclear waste disposal” “the dog crossed the street” และ “Bill Gates” ส่วนตัวอย่างของวลีที่เป็นกริยา เช่น “playing ice hockey” และ “went to school” เป็นต้น

2. วลีทางสถิติ (Statistical Phrase) หมายถึง เซต หรือ ลำดับของคำใดๆ ที่ปรากฏต่อเนื่องกันในเอกสาร โดยมีค่านัยสำคัญทางสถิติ ตัวอย่างวลีประเภทนี้เช่น การสร้าง n-gram ซึ่งคือ ลำดับของคำใดๆ ที่ปรากฏต่อเนื่องกันในเอกสาร เช่น มีเนื้อหาคือ “the dog crossed the street” หลังจากผ่านการเตรียมข้อมูล และตัด Stopwords ออกแล้ว จะได้เนื้อหาเป็น “dog cross street” สร้าง Bi-gram (n=2) จะได้ฟีเจอร์ที่เป็นวลี คือ “dog cross” และ “cross street”

2.4 การเลือกฟีเจอร์ (Feature Selection)

การเลือกฟีเจอร์ มีวัตถุประสงค์คือ ลดจำนวนฟีเจอร์ที่จะใช้แทนเอกสาร เพื่อเพิ่มประสิทธิภาพในการจำแนกประเภทเอกสาร โดยการเลือกฟีเจอร์จะทำโดยเลือกเซตย่อยของฟีเจอร์มาใช้ หรือเรียกได้อีกอย่างหนึ่งว่า Term Space Reduction โดยให้ T คือเซตของเทอมเริ่มต้น การลดฟีเจอร์คือเลือกเซตย่อย T' โดย $|T| \geq |T'|$ [4]

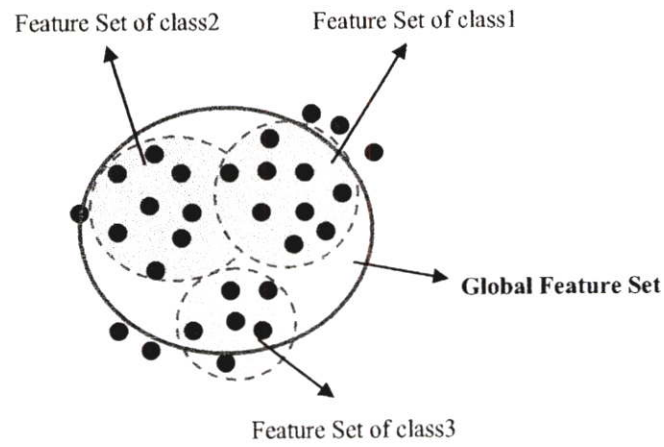
การลดฟีเจอร์มี 2 ระดับ คือ

1. Local Feature Reduction: วิธีนี้จะเลือกเซตย่อยของเทอมภายใต้คลาสต่างๆ โดยเลือกทีละคลาส c_i นั่นคือ ในคลาสต่างๆ จะได้จำนวนฟีเจอร์ไม่เท่ากัน

2. Global Feature Reduction: วิธีนี้จะเลือกเซตย่อยของเทอมภายใต้คลาสทั้งหมด

$$C = \{c_1, \dots, c_l\}$$

รูป 2.3 แสดงการลดฟีเจอร์ทั้ง 2 ระดับ คือ ระดับ Local ที่มีคลาส 3 คลาสด้วยกัน และระดับ Global



รูปที่ 2.3 การเลือกฟีเจอร์ในระดับ Local และ Global

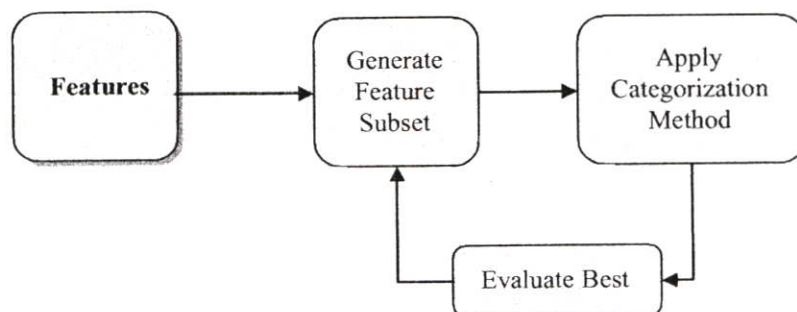
ส่วนวิธีในการเลือกฟีเจอร์นั้นมี 2 แนวทางด้วยกัน คือ

2.4.1 แรพเพอร์ (Wrapper)

ในแนวทางนี้ ชุดย่อยของฟีเจอร์ถูกเลือกโดยดูจากความถูกต้องของตัวจำแนกประเภท และนำอัลกอริทึมที่ใช้ในการจำแนกประเภทจริงๆ มาใช้คิดคำนวณ แนวทางนี้เกิดจากข้อเท็จจริงที่ว่า อัลกอริทึมการเรียนรู้ในการการจำแนกประเภทที่ต่างกัน ได้รับผลกระทบต่างกันจากการเลือกฟีเจอร์ [19]

โดยเริ่มจากเซตของเทอมเริ่มต้น เซตของเทอมใหม่ถูกสร้างขึ้น โดยการเพิ่มหรือลบเทอมออก เมื่อเซตของเทอมใหม่ถูกสร้างขึ้นแล้ว ก็ลองสร้างตัวจำแนกประเภทโดยใช้เซตของเทอมใหม่นี้เป็นฟีเจอร์ แล้วทดสอบความถูกต้อง แล้วเลือกใช้เซตของเทอมที่ให้ประสิทธิภาพดีที่สุดเป็นฟีเจอร์ (Optimal Feature Set) การเลือกเซตย่อยที่ดีที่สุดแสดงดังรูป 2.4

แต่โดยทางเทคนิควิธีการนี้ มีความยากในการอิมพลิเมนต์ โดยเฉพาะอย่างยิ่งเมื่อข้อมูลมีขนาดใหญ่ๆ เช่นมีฟีเจอร์ทั้งหมด m ตัว จะมีเซตย่อยที่เป็นไปได้ 2^m เซต โดยเซตย่อยสามารถถูกค้นหาโดยใช้ Forward Selection หรือ Backward Elimination ซึ่งทำโดยใช้ Hill Climbing หรือใช้ Heuristic อื่นๆ



รูปที่ 2.4 การเลือกเซตย่อยของฟีเจอร์ที่ดีที่สุด [5]

2.4.2 การกรอง (Filter)

ในแนวทางนี้ จะเลือกเซตย่อยของพีเจอร์ โดยทำการกรอง ซึ่งใช้คะแนนเป็นเกณฑ์ ซึ่งคะแนนจะบ่งบอกถึงความสำคัญของเทอม ในการจำแนกประเภทเอกสารมักจะใช้แนวทางนี้ ซึ่งอิมพลีเม้นท์ที่ง่ายกว่าแรพเพอร์ โดยพีเจอร์จะถูกเลือกโดยใช้เกณฑ์ต่างๆ ดังนี้

1. ความถี่ของเอกสาร (Document Frequency) : วิธีนี้พีเจอร์ถูกเลือกโดยดูจากจำนวนเอกสาร โดยเทอมที่ปรากฏในเอกสารส่วนใหญ่จะถูกเก็บไว้ เช่นมีการตั้งค่าจำนวนเอกสารต่ำสุดไว้ แล้วลบเทอมที่ปรากฏแค่เพียงตามจำนวนเอกสารที่ตั้งไว้ ออก

2. ใช้ตัววัดจาก Information-theoretic Function (f) เช่น

Information Gain Measure: ค่า Information Gain $IG(t, c_i)$ ของเทอม t ในคลาส c_i หาได้ดังนี้ [8]

$$IG(t, c_i) = \sum_{c \in \{c_i, c_j\}} \sum_{t \in \{t_k, t_l\}} P(t, c_i) \cdot \log \frac{P(t, c_i)}{P(t) \cdot P(c_i)} \quad (2.1)$$

Mutual Information Measure: ค่า Mutual Information $MI(t, c_i)$ ของเทอม t ในคลาส c_i หาได้จากสมการต่อไปนี้ [9]

$$MI(t, c_i) = \log \frac{P(t, c_i)}{P(t) \cdot P(c_i)} \quad (2.2)$$

Odds Ratio Measure: ค่า Odds Ratio $OR(t, c_i)$ ของเทอม t ในคลาส c_i หาได้จากสมการต่อไปนี้ [10]

$$OR(t, c_i) = \frac{P(t | c_i) \cdot (1 - P(t | \bar{c}_i))}{(1 - p(t | c_i)) \cdot p(t | c_i)} \quad (2.3)$$

โดย $P(t, c_i)$ คือ ความน่าจะเป็นที่คลาส c_i และเทอม t เกิดขึ้นร่วมกัน

$P(t | c_i)$ คือ ความน่าจะเป็นที่เทอม t จะเกิดเมื่อเกิดคลาส c_i

$P(t | \bar{c}_i)$ คือ ความน่าจะเป็นที่เทอม t จะเกิดเมื่อไม่เกิดคลาส c_i

$P(t)$ คือ ความน่าจะเป็นในการเกิดเทอม t

$P(c_i)$ คือ ความน่าจะเป็นของคลาส c_i

ซึ่งฟังก์ชันตัววัดที่กล่าวมาจะเป็นแบบ Local ซึ่งเราสามารถทำแบบ Global ได้โดยนำค่าที่ได้จากแต่ละคลาสมารวมกัน ซึ่งอาจเพิ่มการถ่วงน้ำหนักให้แต่ละคลาส หรือใช้ค่าสูงสุดก็ได้ ดังนี้

$$f_{sum}(t) = \sum_{i=1}^{|C|} f(t, c_i) \quad (2.4)$$

$$f_{wsum}(t) = \sum_{i=1}^{|C|} P(c_i) f(t, c_i) \quad (2.5)$$

$$f_{max}(t) = \max_{i=1}^{|C|} f(t, c_i) \quad (2.6)$$

โดย $f(t, c_i)$ คือ ฟังก์ชันของตัววัดที่ใช้คำนวณค่าของเทอม t ในคลาส c_i
 $f_{sum}(t)$ คือ ค่าที่ได้จากการรวมค่าจากฟังก์ชัน f ของเทอม t จากทุกคลาส
 $f_{wsum}(t)$ คือ ค่าที่ได้จากการรวมค่าจากฟังก์ชัน f ของเทอม t จากทุกคลาสแบบถ่วงน้ำหนัก
 $f_{max}(t)$ คือ ค่าสูงสุดจากฟังก์ชัน f ของเทอม t จากทุกคลาส

2.5 การแทนเอกสาร (Document Representation)

จากพีเจอร์ที่เราได้จากกระบวนการที่ผ่านมา ต่อไปเราจะทำการแทนเอกสารโดยใช้พีเจอร์เหล่านี้ โดยทั่วไปวิธีที่ใช้กันอย่างแพร่หลายคือวิธี Vector Space Model [2] ในโมเดลนี้แต่ละเอกสารจะถูกแทนโดยเวกเตอร์ d โดย d เขียนได้ดังนี้

$$d = \{w_1, w_2, \dots, w_m\}$$

โดย w_i ($i=1, 2, \dots, m$) คือ ค่าน้ำหนักเทอมของเทอม t_i ที่เป็นพีเจอร์ และ m คือ จำนวนพีเจอร์ แต่ละค่าถ่วงน้ำหนักเทอมแทนนัยสำคัญของแต่ละเทอมในเอกสาร ซึ่งในการคำนวณค่าถ่วงน้ำหนักเทอมนี้จะพิจารณาความถี่ในการปรากฏของเทอม ภายในแต่ละเอกสาร และในเซตของเอกสารทั้งหมด วิธีการต่างๆ ในการกำหนดค่าน้ำหนัก มีดังนี้

2.5.1 Boolean Weighting

เป็นวิธีการที่ง่ายที่สุด วิธีนี้สามารถบอกได้เพียงว่ามีพีเจอร์นั้นปรากฏในเอกสารหรือไม่ คือแทนค่า Weight ด้วย 0 และ 1 โดยจะแทนค่าเป็น 1 ถ้ามีพีเจอร์นั้นปรากฏในเอกสาร และแทนค่าเป็น 0 เมื่อไม่มีพีเจอร์นั้นปรากฏในเอกสาร ดังสมการที่ (2.7)

$$w_{ij} = \begin{cases} 1 & \text{if } tf_{ij} > 0 \\ 0 & \text{otherwise} \end{cases} \quad (2.7)$$

โดย w_{ij} คือ ค่าน้ำหนักของเทอม i ในเอกสาร j และ
 tf_{ij} คือ ความถี่เทอม i ในเอกสาร j

2.5.2 Term Frequency Weighting

วิธีการนี้สามารถบอกความถี่ของพีเจอร์ที่ปรากฏในเอกสารได้ จึงมีประสิทธิภาพดีกว่าแบบ Boolean Weighting เป็นวิธีการที่หาความถี่ของพีเจอร์ที่ปรากฏในเอกสารว่ามีพีเจอร์นั้นปรากฏในเอกสารกี่ครั้ง โดยยังคงใช้ 0 แทนความหมายว่าไม่มีพีเจอร์นั้นปรากฏในเอกสาร ดังสมการที่ (2.8)

$$w_{ij} = tf_{ij} \quad (2.8)$$

2.5.3 TF-IDF Weighting

วิธีนี้เป็นที่นิยมมากที่สุด โดยการใช้ความถี่ของเทอม (Term Frequency) ร่วมกับความถี่เอกสารผกผัน (Inverse Document Frequency) หรือเรียกย่อๆ ว่า TF-IDF [11] ดังสมการที่ (2.9)

$$w_{ij} = tf_{ij} * \log\left(\frac{N}{n_i}\right) \quad (2.9)$$

โดย N คือจำนวนเอกสารทั้งหมด และ n_i คือ จำนวนเอกสารที่มีเทอม i ปรากฏอยู่

2.5.4 TFC Weighting

วิธีการนี้จะคล้ายกับ TF-IDF Weighting แต่มีการ Normalization เพราะว่าเอกสารมีขนาดไม่เท่ากัน [11] ดังสมการที่ (2.10)

$$w_{ij} = \frac{tf_{ij} * \log\left(\frac{N}{n_i}\right)}{\sqrt{\sum_{k=1}^M \left[tf_{kj} * \log\left(\frac{N}{n_k}\right)\right]^2}} \quad (2.10)$$

วิธีการทั้งหมดที่กล่าวมาเป็นวิธีการในการกำหนดค่าน้ำหนักของเทอมที่ใช้เป็นพีเจอร์ ที่เป็นตัวแทนของเอกสาร

ตัวอย่างการคำนวณค่าน้ำหนัก มีดังนี้ มีเอกสารอยู่ 3 เอกสาร ทำการเตรียมเอกสารแล้วได้เนื้อหา และเทอมดังในตารางที่ 2.1

ตารางที่ 2.1 แสดงตัวอย่างเอกสาร เนื้อหาและเทอมที่ได้หลังการเตรียมเอกสาร

document	text	terms
d1	web web graph	web graph
d2	graph web net graph net	graph web net
d3	page web complex	page web complex

จากตัวอย่างในตาราง 2.1 เทอมที่ได้จากแต่ละเอกสารจะถูกนำมาใช้เป็นฟีเจอร์ สมมติในที่นี้ไม่ได้ทำการเลือกฟีเจอร์ คือ ใช้ทุกเทอมที่ได้ ดังนั้นเราจะได้เซตของฟีเจอร์ ดังนี้

$$Features = \{“web”, “graph”, “net”, “page”, “complex”\}$$

ในการให้ค่าน้ำหนักแบบบูลีน เราจะได้เวกเตอร์ของเอกสาร ดังนี้

$$d_1 = \{1, 1, 0, 0, 0\}$$

$$d_2 = \{1, 1, 1, 0, 0\}$$

$$d_3 = \{1, 0, 0, 1, 1\}$$

ส่วนในการให้ค่าน้ำหนักโดยใช้ความถี่ เราจะได้เวกเตอร์ของเอกสารเป็น

$$d_1 = \{2, 1, 0, 0, 0\}$$

$$d_2 = \{1, 2, 2, 0, 0\}$$

$$d_3 = \{1, 0, 0, 1, 1\}$$

ในการให้ค่าน้ำหนักโดยใช้ TF-IDF เราจะได้เวกเตอร์ของเอกสาร ดังนี้

$$d_1 = \{0, 0.4055, 0, 0, 0\}$$

$$d_2 = \{0, 0.8109, 2.1972, 0, 0\}$$

$$d_3 = \{0, 0, 0, 1.0986, 1.0986\}$$

ซึ่งจะสังเกตได้ว่าคำที่ปรากฏทุกเอกสารจะถูกให้ค่าน้ำหนักเท่ากับศูนย์ เพราะไม่สามารถใช้แบ่งแยกเอกสารได้

และในการให้ค่าน้ำหนักโดยใช้ TFC เราจะได้เวกเตอร์ของเอกสารเป็น

$$d_1 = \{0, 1, 0, 0, 0\}$$

$$d_2 = \{0, 0.3462, 0.9381, 0, 0\}$$

$$d_3 = \{0, 0, 0, 0.7071, 0.7071\}$$

2.6 วิธีการจำแนกประเภทเอกสาร (Categorization Method)

หลายๆ วิธีการในการจำแนกประเภทได้ถูกนำมาใช้ในการจำแนกประเภทเอกสาร เช่น Naïve Bayes (NB) Classifier, Decision Tree Classifier, Neural Network และ k-Nearest Neighbor Classifier เป็นต้น

ในงานวิจัยนี้เราได้ใช้ k-Nearest Neighbor Classifier ในการจำแนกประเภทเอกสาร ซึ่งมีรายละเอียดดังต่อไปนี้

การจำแนกประเภทแบบ k-Nearest Neighbor (k-NN) [5] เป็นวิธีการเรียนรู้จากกรณีตัวอย่าง (Instance-based Learning Method) ซึ่งมีการเก็บชุดข้อมูลที่จะใช้เอาไว้ก่อน จนกว่าจะมี

ข้อมูลใหม่เข้ามา จึงจะเริ่มกระบวนการจำแนกประเภท ทำให้มันถูกเรียกว่าเป็นการเรียนรู้แบบขี้เกียจ (Lazy Learning) การจำแนกประเภทแบบ k-NN นี้ ประสิทธิภาพจะขึ้นอยู่กับตัววัดความคล้ายคลึงที่ใช้ (Similarity Measure)

ขั้นตอนแรกในการจำแนกโดยวิธี k-NN มัน จะคำนวณความคล้ายคลึงระหว่างเอกสารใหม่กับเอกสารทุกตัวที่เราทราบว่าอยู่ประเภทไหนแล้ว และเลือกเอกสารที่มีค่าความคล้ายคลึงสูงที่สุดมา k ตัว จากนั้นก็จะทำการเลือกประเภทจาก k เอกสารนั้นให้กับเอกสารใหม่ ซึ่งกระบวนการจำแนกประเภทเอกสารแบบ k-NN เขียนเป็นอัลกอริทึมได้ดังรูป 2.5

Algorithm: kNN Text Categorization

Input: D = Vectors of documents, k, new document

Output: Class of new document

For all Document in D

Compute Similarity between

each document and new document

End

Select the k documents that are most similar to new document

Select class of new document from that k documents

รูปที่ 2.5 อัลกอริทึมการจำแนกประเภทเอกสารแบบ kNN

ซึ่งในการเลือกประเภทเอกสารให้เอกสารใหม่จากทั้ง k เอกสาร (kNN) นั้น มี 2 วิธี ที่ใช้กันทั่วไป ดัง 2 สมการต่อไปนี้

$$C(d_i) = \arg \max_{i=1}^{|C|} \sum_{x_j \in kNN} y(x_j, c_i) \quad (2.13)$$

$$C(d_i) = \arg \max_{i=1}^{|C|} \sum_{x_j \in kNN} \text{sim}(d_i, x_j) \cdot y(x_j, c_i) \quad (2.14)$$

โดย d_i คือ เอกสารใหม่ที่จะทำการจำแนกประเภท

x_j คือ หนึ่งในเอกสารใกล้เคียงที่เลือกมา k ตัว

$\text{sim}(d_i, x_j)$ คือ ความคล้ายคลึงระหว่างเอกสาร d_i กับเอกสาร x_j

$y(x_j, c_i) \in \{0,1\}$ หมายถึงว่าเอกสาร x_j นั้นอยู่ในคลาส c_i หรือไม่ ถ้ามีค่าเท่ากับ 1 แสดงว่าเอกสาร x_j นั้นอยู่ในคลาส c_i แต่ถ้ามีค่าเท่ากับ 0 แสดงว่าเอกสาร x_j ไม่ได้อยู่ในคลาส c_i

สมการที่ (2.13) หมายถึงว่า เอกสารใหม่นั้นจะถูกจำแนกประเภทให้อยู่ในคลาสที่มีจำนวนสมาชิกที่อยู่ใน เอกสารใกล้เคียง k เอกสารนั้นมากที่สุด ในขณะที่สมการที่ (2.14) นั้น หมายถึงว่า เอกสารใหม่นั้นจะถูกจำแนกให้อยู่ในคลาสที่มีผลรวมค่าความคล้ายคลึงของคลาสในเอกสาร ใกล้เคียง k เอกสารมากที่สุด

จะเห็นได้ว่าการวัดความคล้ายคลึงของเอกสารเป็นส่วนสำคัญที่เราใช้ในการจำแนกประเภทเอกสาร การวัดความคล้ายคลึงของเอกสารที่นิยมใช้ มี 2 วิธีหลักๆ คือ

1. การวัดความคล้ายคลึง โดยใช้ระยะทาง Euclidean

ให้ขนาดของเวกเตอร์ คือจำนวนเทอมเท่ากับ m ความคล้ายคลึงระหว่างเวกเตอร์เอกสาร d_1 และ d_2 หรือ $\text{sim}(d_1, d_2)$ โดยใช้การวัดระยะทาง Euclidean [12] มีสมการดังนี้

$$\text{sim}(d_1, d_2) = \sqrt{\frac{\sum_{i=1}^m (w_{i1} - w_{i2})^2}{m}} \quad (2.11)$$

2. การวัดความคล้ายคลึง โดยใช้สหสัมพันธ์โคไซน์

การวัดสหสัมพันธ์โคไซน์ (Cosine Correlation Measure) [13] เป็นอีกวิธีหนึ่งที่เป็นที่นิยมใช้ในการวัดความคล้ายคลึงของเอกสาร เขียนเป็นสมการได้ดังนี้

$$\text{sim}(d_1, d_2) = \frac{\sum_{i=1}^m w_{i1} \cdot w_{i2}}{\sqrt{\sum_i w_{i1}^2} \cdot \sqrt{\sum_i w_{i2}^2}} \quad (2.12)$$

โดย $\text{sim}(d_1, d_2)$ คือ ความคล้ายคลึงระหว่างเอกสาร d_1 กับเอกสาร d_2

w_{i1} และ w_{i2} คือ ค่าน้ำหนักของเทอมที่ i ในเอกสาร d_1 กับเอกสาร d_2 ตามลำดับ

m คือ จำนวนฟีเจอร์ที่ใช้แทนเอกสาร

2.7 การจำแนกประเภทเอกสารโดยใช้วลี

การจำแนกประเภทเอกสาร โดยใช้วลี มีวัตถุประสงค์เพื่อแก้ปัญหาที่เกิดจากการจำแนกประเภทเอกสาร โดยใช้คำ ข้อเสียของการแทนเอกสารโดยการใช้คำ คือ การทำลายความสัมพันธ์เชิงความหมายระหว่างคำไป โดยวลีที่มีความหมายแน่นอน เช่น “White House” หรือ “Bill Gates” เมื่อถูกแทนโดยการใช้คำ วลีเหล่านี้ก็จะถูกแยกคำออก ทำให้ความหมายของพวกมันเสียไป ข้อเสีย

ของการใช้คำอีกประการหนึ่งคำคำหนึ่งอาจมีได้หลายความหมาย ซึ่งถ้าเป็นวลีจะบอกถึงความหมายที่แท้จริงได้มากกว่า เช่น คำว่า “Java” เป็นไปได้ 2 ความหมาย คือ ภาษาโปรแกรมมิ่งหรือ เกาะของประเทศอินโดนีเซีย โดยถ้าเราใช้วลี “Java Script” จะบอกความหมายที่แท้จริงคือในความหมายแรกได้

วลีที่ใช้ในการแทนเอกสาร แบ่งได้ 2 ลักษณะ คือ วลีทางภาษา (Syntactic Phrase) และ วลีทางสถิติ (Statistical Phrase) โดยในการใช้พีเจอร่วมนั้น ยังมี 2 แนวทางด้วยกัน คือ

1. ใช้พีเจอรที่เป็นวลีเพียงอย่างเดียวในการแทนเอกสาร
2. ใช้พีเจอรที่เป็นวลีร่วมกับพีเจอรที่เป็นคำในการแทนเอกสาร เนื่องจาก จากงานวิจัยที่ศึกษาการใช้ Syntactic Phrase ในการแทนเอกสาร เพียงอย่างเดียว ผลที่ได้กลับทำให้ความถูกต้องในการจำแนกประเภทเอกสารน้อยกว่าการใช้คำเป็นพีเจอรเพียงอย่างเดียว งานวิจัยเรื่อง “An Evaluation of Phrasal and Clustered Representations on a Text Categorization Task” ของ Lewis [24] ใช้วลีที่เป็นคำนามที่ปรากฏอย่างน้อย 2 ครั้ง ในการแทนเอกสาร สำหรับการจำแนกประเภทเอกสาร แต่ผลที่ได้กลับแย่ลงกว่าการใช้คำเป็นพีเจอรเพียงอย่างเดียว เขาได้ให้เหตุผลไว้ว่า

- ก) เทอมที่ได้มีปริมาณมาก
- ข) ค่าของพีเจอรกระจายไม่เท่ากัน โดยหลายๆ คำมีค่าต่ำมาก
- ค) พีเจอรมีความซ้ำซ้อนกันสูง
- ง) พีเจอรประกอบด้วย Noise

Lewis ได้พยายามแก้ปัญหาเหล่านี้โดยใช้อัลกอริทึมการจัดกลุ่ม (Clustering) ในการรวมพีเจอรที่เป็นวลี ไปเป็น Meta-features แต่ในการแทนเอกสารแบบใหม่นี้ก็ไม่ได้สร้างผลที่ดีกว่าเท่าใดนัก

ในการศึกษาเปรียบเทียบระหว่างการใช้วลีทางภาษากับวลีทางสถิติ โดย Mitra และกลุ่ม ได้ทำการทดลองเปรียบเทียบการใช้ระหว่าง Bi-grams กับวลีที่เป็น POS-tag Sequence (เช่น Noun-Noun, Adjective-Noun) ในงาน Information Retrieval ผลที่ได้คือ การใช้วลีทางภาษาดีกว่าเพียงเล็กน้อย ซึ่งสรุปแล้วผลที่ได้น่าจะขึ้นอยู่กับลักษณะการใช้งานและข้อมูลที่เหมาะสมกับวลีแบบไหน แต่จากผลของ Lewis ที่แสดงให้เห็นว่า การใช้วลีเพียงอย่างเดียวนั้น ไม่ได้เพิ่มประสิทธิภาพในการจำแนกเอกสาร จึงมีหลายงานวิจัยที่ได้แนะนำว่าควรที่จะเพิ่มความถูกต้องของการจำแนกประเภทเอกสาร โดยการใช้วลีทางสถิติร่วมกับการใช้คำเดียวเป็นพีเจอร เช่นงานวิจัยของ Mladeni และ Grobelink รายงานการวิจัยเรื่อง “Word Sequences as Features in Text Learning” [25] ได้สร้างพีเจอรใหม่จากลำดับของคำ ตั้งแต่ 2-5 คำ (n-gram, n=5) โดยเลือกใช้เฉพาะพีเจอรที่มีคะแนนสูง โดยในที่นี้คะแนนที่ใช้ คือ ความถี่ของเทอม ข้อมูลที่ใช้คือ Yahoo Text Hierarchy และใช้ Naïve Bayes Classifier ในการจำแนกประเภทเอกสาร พวกเขาแสดงให้เห็นว่าการใช้ลำดับของ

สำนักหอสมุดกลาง พระจอมเกล้าลาดกระบัง

คำถึง 3 (n-gram, n=3) ร่วมกับคำเดี่ยว แทนที่การใช้คำเดี่ยวเพียงอย่างเดียว นั้น สามารถเพิ่มประสิทธิภาพในการจำแนกประเภทได้ โดยลำดับที่ยาวกว่า โดยเฉลี่ยแล้วไม่ส่งผลกระทบต่อประสิทธิภาพในการจำแนกประเภท

2.7.1 การลดพีเจอร์วลีสำหรับการจำแนกประเภทเอกสาร

เนื่องจากการใช้พีเจอร์วลีที่เป็นวลียังมีปัญหาในการนำมาใช้คือ เซตของพีเจอร์วลีมีขนาดใหญ่ ซึ่งมีพีเจอร์วลีที่มีความสามารถในการจำแนกประเภทต่ำอยู่มาก ดังนั้นจึงได้มีหลายงานวิจัยที่ได้พยายามสร้างและเลือกพีเจอร์วลี ที่จะใช้ร่วมกับคำเดี่ยวในการแทนเอกสาร ให้มีจำนวนน้อยลง และมีประสิทธิภาพเพิ่มขึ้น เช่น

งานวิจัยเรื่อง “*The Use of Bi-grams to Enhance Text Categorization*” โดย Chade-Meng Tan และกลุ่ม [26] ได้ทำการเลือกคำเดี่ยว โดยเลือกใช้คำเดี่ยวที่มีความถี่ของเอกสารสูงเป็นอันดับต้นๆ และใช้คำเดี่ยวที่เลือกนั้น มาสร้าง Bi-gram จากคำที่ติดกันในเอกสารที่ผ่านการเตรียมและลบ Stopwords ออกแล้ว โดยคำใดคำหนึ่งจะต้องอยู่ภายในเซตของคำเดี่ยวที่เลือกไว้ และทำการเลือก Bi-grams โดยใช้เกณฑ์ความถี่และค่า Information Gain ในแต่ละคลาส โดยมีการตั้งค่าต่ำสุดของแต่ละค่าไว้ ชุดข้อมูลที่ใช้ คือ เว็บเพจ จาก Yahoo! Science Hierarchy และข้อมูล Reuters-21578 ผลการทดลองจำแนกประเภทเอกสารโดยใช้ Naïve Bayes แสดงให้เห็นว่า วลีที่ถูกเลือกอย่างเหมาะสมรวมกับคำเดี่ยว ให้ผลดีกว่าการใช้คำเดี่ยวเพียงอย่างเดียวเป็นพีเจอร์

ต่อมา Bekkerman และ Allan ได้รายงานเรื่อง “*Using Bigrams in Text Categorization*” [27] โดยทำการสร้าง Bi-grams จากคำเดี่ยวที่อยู่ในลำดับต้นๆ โดยดูจากค่า Mutual Information และเลือกใช้เฉพาะ Bi-gram ที่ดีกว่าคำที่เป็นส่วนประกอบของมันทั้ง 2 คำ โดยดูจากค่า Mutual Information และเลือก Bi-grams ที่ได้ในอันดับต้นๆ มาใช้

ชุดข้อมูลที่ใช้ในการทดลอง คือ 20 Newsgroups โดยเลือกใช้คำเดี่ยว 5000 คำ และ เลือก Bi-grams 1000 คำ ที่ได้จากการเลือกจากแต่ละคลาสมาใช้เป็นพีเจอร์วลีร่วมกัน ผลจากการทดลองโดยใช้ SVM Classifier แสดงให้เห็นว่าการใช้วลีร่วมกับคำให้ผลการจำแนกประเภทที่ดีกว่า การใช้คำเพียงอย่างเดียว

จากงานวิจัยที่ผ่านมา เราสามารถสรุปปัญหาในการใช้วลีเป็นพีเจอร์วลีในการแทนเอกสารได้ ดังนี้

1. วลีมีความจำเพาะเจาะจงเกินไป คือ มีความหมายเฉพาะมากเกินไป เช่น วลีทางภาษา หรือวลีที่ได้จากเอกสารโดยตรง ทำให้คำที่ถูกแทนในเอกสารส่วนใหญ่ มีค่าเป็นศูนย์ ตัวอย่างเช่น วลี “*k-nearest neighbor categorization algorithm*” จะปรากฏในแต่ละเอกสารน้อยมาก แต่ถ้าเราใช้วลีย่อยเป็น “*k-nearest neighbor*” หรือ “*categorization algorithm*” จะมีการกระจายของค่าในการ

แทนเอกสารดีกว่า เพราะมีโอกาสในการปรากฏมากกว่าในแต่ละเอกสาร และในการสกัดวลีทางภาษายังทำได้ยากด้วย

2. แต่ถ้าวลีนั้นสั้นไป ก็จะเกิดปัญหาคล้ายกับการใช้คำเดี่ยว คือ มีความสามารถในการแบ่งแยกประเภทต่ำ โดยเฉพาะอย่างยิ่งเมื่อประเภทของข้อมูลมีความใกล้เคียงกันมาก เช่น วลี “*drama player voting*” อยู่ในข่าวบันเทิง ส่วนวลี “*football player voting*” อยู่ในข่าวกีฬา ถ้าเราใช้ฟิเจอร์เป็น “*player voting*” ก็ทำให้ประสิทธิภาพในการจำแนกว่าอยู่ในข่าวประเภทไหนลดลง

3. จากปัญหาข้อ 1 กับ 2 ถ้าวลียาวไปก็ไม่ได้ สั้นเกินไปก็ไม่ได้ ทางแก้หนึ่งก็คือ ใช้ทั้งวลี และวลีย่อยเป็นฟิเจอร์ แต่ก็เกิดปัญหาความซ้ำซ้อนของวลีขึ้นมาอีก ดังนั้น จึงต้องมีการเลือกว่า ควรจะใช้วลีไหนดีในการแทนเอกสาร หรือใช้ทั้ง 3 ตัวอย่างในกรณีตัวอย่างข้อที่ 1 คือใช้ ทั้งวลี “*k-nearest neighbor categorization algorithm*” “*k-nearest neighbor*” และ “*categorization algorithm*” แทนเอกสาร

4. การใช้วลีทางสถิติเป็นการแก้ปัญหาอย่างหนึ่งของการใช้วลีทางภาษา เพราะไม่ต้องคำนึงถึงหลักไวยากรณ์ทางภาษา เช่น การสร้าง n -gram ซึ่งคือ ลำดับของคำใดๆ ที่ปรากฏต่อเนื่องกันในเอกสาร โดยส่วนใหญ่จะใช้ $n=2$ เรียกว่า Bi-gram แต่ปัญหาก็คือ วลีที่สร้างได้มีจำนวนมาก ที่เป็นวลีขยะ คือ วลีที่มีอำนาจในการจำแนกประเภทต่ำ และในการสร้าง n -gram นั้น จำเป็นที่จะต้องกำหนดขนาดสูงสุดของวลีไว้ คือ n ทำให้บางครั้งมีวลีที่มีความสามารถในการจำแนกประเภทที่ดี แต่มีขนาดมากกว่า n ก็จะไม่ถูกใช้

ดังนั้น พอจะสรุปได้ว่าฟิเจอร์ที่จะนำมาใช้แทนเอกสารนั้น ควรจะมีคุณสมบัติ 2 ประการคือ

1. มีอำนาจในการแบ่งแยก (Discriminant Power) สูง คือ วลีนั้นจะต้องไม่มีความหมายแคบเกิน ทำให้มีความถี่น้อย และจะต้องไม่มีความหมายกว้างจนเกินไป
2. มีปริมาณน้อย โดยไม่กระทบต่อความถูกต้องในการจำแนกประเภทเอกสาร ซึ่งต้องอาศัยการเลือกฟิเจอร์ที่มีประสิทธิภาพ

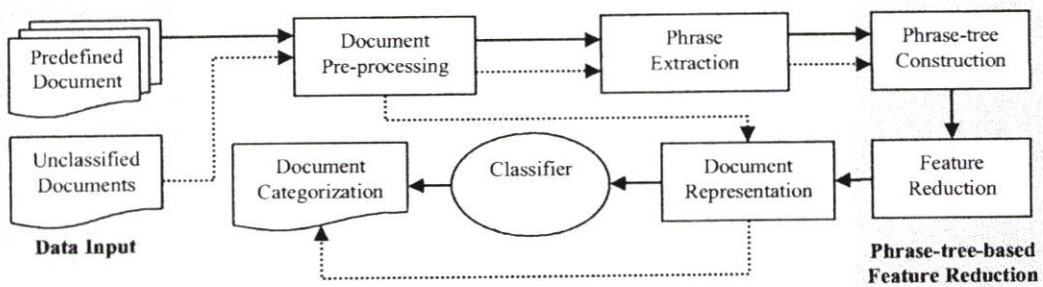
ในงานวิจัยนี้พยายามแก้ปัญหาดังกล่าว โดยนำ ทฤษฎีต่างๆ ที่เกี่ยวข้อง เหล่านี้ มาใช้ในการแก้ปัญหาในการจำแนกประเภทเอกสาร โดยทำการลดฟิเจอร์โดยใช้ต้นไม้วลีสำหรับการจำแนกประเภทเอกสาร ซึ่งจะกล่าวถึงรายละเอียดในบทที่ 3 ต่อไป

บทที่ 3

การลดพีเจอร์โดยใช้ต้นไม้วลีสำหรับการจำแนกประเภทเอกสาร

ในงานวิจัยนี้ได้พยายามแก้ไขปัญหามาจากการใช้วลีเป็นพีเจอร์ในการแทนเอกสาร โดยนำเสนอวิธีการสร้างพีเจอร์วลีแบบใหม่ ที่ได้ปริมาณน้อย แต่ยังคงมีประสิทธิภาพในการจำแนกประเภทเอกสาร คือ มีอำนาจในการแบ่งแยกสูง โดยมีส่วนที่สำคัญ คือ เทคนิคในการสกัดวลี การสร้างพีเจอร์และการเลือกวลีโดยใช้โครงสร้างข้อมูลแบบต้นไม้

กระบวนการของการจำแนกประเภทเอกสาร โดยทำการลดพีเจอร์โดยใช้ต้นไม้วลี แสดงดังรูป 3.1 โดยเส้นทึบแทนการไหลของข้อมูลในขั้นตอนการสร้างตัวจำแนกประเภท โดยเอกสารที่เรารู้คลาสแล้วจะผ่านกระบวนการเตรียมเอกสาร จากนั้นผ่านกระบวนการการลดพีเจอร์โดยใช้ต้นไม้วลี และการแทนเอกสาร เพื่อสร้างเป็นตัวจำแนกประเภท ส่วนเส้นประหมายถึงเส้นทางการไหลของข้อมูลในขั้นตอนการจำแนกประเภท โดยเมื่อมีเอกสารที่เรายังไม่ทราบคลาสเข้ามา ก็ผ่านกระบวนการเตรียมเอกสาร และการแทนเอกสาร เพื่อนำเข้าสู่ตัวจำแนกประเภท ทำให้ทราบประเภทของเอกสารนั้น



รูปที่ 3.1 กระบวนการของการจำแนกประเภทเอกสาร โดยทำการลดพีเจอร์โดยใช้ต้นไม้วลี

3.1 การเตรียมเอกสาร (Data Pre-processing)

ในกระบวนการนี้ เอกสารจะถูกเตรียมและแปลงให้อยู่ในรูปแบบที่จะเข้าสู่กระบวนการการสกัดวลีต่อไป โดยมีรายละเอียดดังต่อไปนี้

1. ในกรณีที่เอกสารเป็นเว็บลบบส่วนที่เป็นแท็ก (Tag) ของเอกสารที่เป็นเว็บออกก่อน เช่น “<html>”, “</html>”, “
”, “”
2. แปลงให้อักษรตัวพิมพ์เล็กทั้งหมด และลบตัวเลขออก
3. เปลี่ยนเครื่องหมายต่างๆ เช่น , () [] { } / \ | = ‘ ’ “ ” ^ ? < > และ Stopwords ให้เป็นเครื่องหมาย . (Full Stop Point) เพียงอย่างเดียวเพื่อใช้แบ่งแยกวลีต่อไป แต่จะมีกรณีพิเศษคือเครื่องหมาย - _ & + * ~ # \$ % ซึ่งจะไม่ถูกเปลี่ยน แต่จะถูกลบเครื่องหมายออก เช่น “black-

market” → “black market”, “black&white game” → “black white game”, “bottle + milk = bottle of milk” → “bottle milk . bottle milk” การที่เราไม่ใช้เครื่องหมาย - & + มาแบ่งแยกวลี เพราะทั้งหมดเป็นการบ่งบอก ถึงการเกิดขึ้นร่วมกันของคำ ซึ่งเป็นวลีที่ส่วนมากได้ความหมายอยู่แล้ว ดังนั้น จึงยังไม่ต้องแยกออก

ในการแบ่งแยกวลีแบบนี้ คล้ายกับงานของ Klok ที่ใช้การสกัดวลีโดยใช้เทอมที่อยู่ระหว่าง Stopwords มาเป็นวลี ที่ใช้ในการสืบค้นร่วมกับคำและวลีที่เป็นคำนาม [21]

4. ทำแต่ละคำให้อยู่ในรากศัพท์ (Stemming) เพื่อลดจำนวนเทอมที่ได้

3.2 การสกัดวลี (Phrase Extraction)

ในกระบวนการการสกัดวลีนี้ วลีจะถูกสกัดจากเอกสารที่ถูกรวบรวมแล้ว โดยมีเครื่องหมาย “.” เป็นตัวแบ่งแยกวลี

สมมติมีตัวอย่างเอกสารดังนี้

“<html> 1. k-Nearest neighbor is an efficient categorization algorithm. </html>”

หลังจากการทำการเตรียมเอกสาร เราจะได้เอกสารที่เตรียมแล้ว คือ

“k-nearest neighbor . effici categor algorithm.”

เมื่อเข้าสู่กระบวนการสกัดวลีนี้ เราจะได้ 2 วลี คือ “k-nearest neighbor” และ “effici categor algorithm” ซึ่งในการสกัดวลีแบบนี้ เราจะได้คำเดี่ยวมาด้วย เช่น เอกสารเป็น

“This is categorization method (classification)”

หลังจากการทำการเตรียมเอกสาร เราจะได้เอกสารที่เตรียมแล้ว คือ

“categor method .classifi.”

เมื่อเข้าสู่กระบวนการสกัดวลีนี้ เราจะได้ทั้งวลี และคำเดี่ยว คือ “categor method” และ “classifi” โดยในที่นี้เราจะรวมคำเดี่ยวที่สกัดได้จากกระบวนการนี้เข้าไปด้วย เพื่อใช้ในการสร้างต้นไม้วลีต่อไป

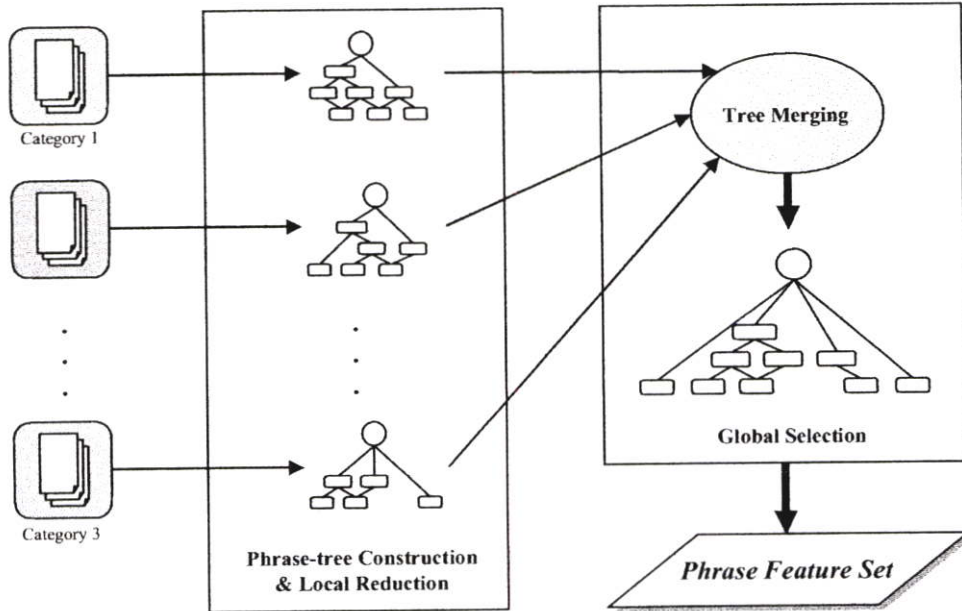
3.3 การลดพีเจอร์โดยใช้ต้นไม้วลีสำหรับการจำแนกประเภทเอกสาร

(Phrase-tree-based Feature Reduction for Text Categorization)

ในการสร้างต้นไม้วลีนั้น มีวัตถุประสงค์เพื่อใช้โครงสร้างข้อมูลแบบต้นไม้ มาช่วยในการสร้างพีเจอร์ โดยโครงสร้างต้นไม้สามารถแสดงถึงความสัมพันธ์ระหว่างวลีและวลีย่อย และคำเดี่ยว ซึ่งช่วยในการเลือกพีเจอร์ที่ลดความซ้ำซ้อนของวลีได้ โดยเราต้องการพีเจอร์ที่เป็นตัวแทน

เอกสารที่ดี มีอำนาจในการจำแนกประเภท โดยจะทำการเลือกฟีเจอร์ทั้งในระดับ Local และระดับ Global โดย

ในการทำการลดฟีเจอร์โดยใช้ต้นไม้วลีสำหรับการจำแนกประเภทเอกสาร ประกอบด้วย 3 ขั้นตอนหลักๆ คือ 1. การสร้างต้นไม้วลี 2. การลดฟีเจอร์ในระดับ Local 3. การลดฟีเจอร์ในระดับ Global แสดงดังรูป 3.2



รูปที่ 3.2 ขั้นตอนการลดฟีเจอร์โดยใช้ต้นไม้วลี

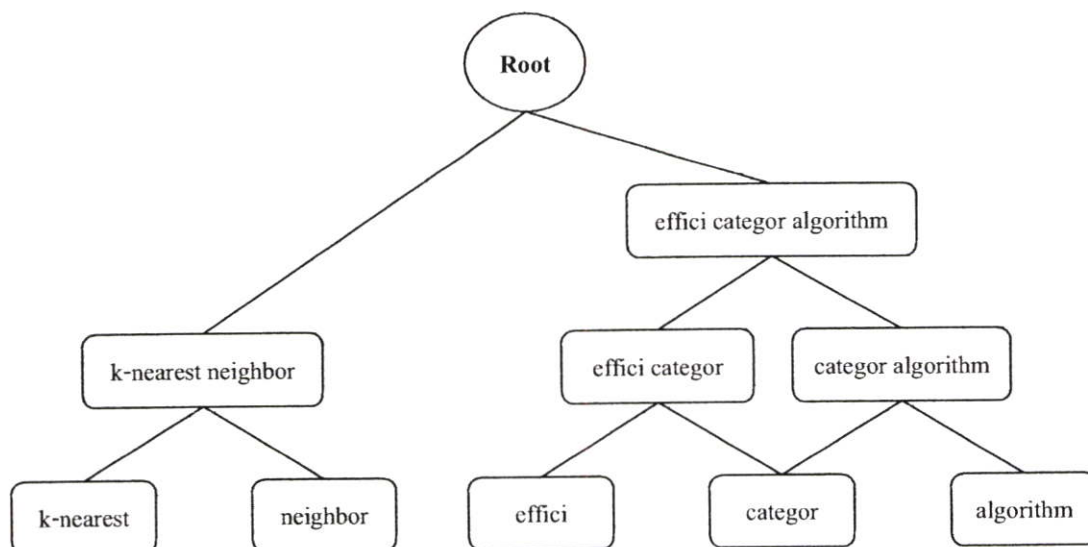
โดยในการทำการลดฟีเจอร์โดยใช้ต้นไม้วลีสำหรับการจำแนกประเภทเอกสาร สามารถเขียนเป็นอัลกอริทึมได้ดังรูปที่ 3.9 ประกอบด้วย 3 ขั้นตอนคือ 1. การสร้างต้นไม้วลี 2. การลดฟีเจอร์ในระดับ Local 3. การลดฟีเจอร์ในระดับ Global และมีรายละเอียดในแต่ละขั้นตอน ดังนี้

3.3.1 การสร้างต้นไม้วลี (Phrase-tree Construction)

ในการสร้างต้นไม้วลี จะเป็นการสร้างต้นไม้วลีของแต่ละคลาส โดยสร้างจากวลีที่สกัดได้จากทุกเอกสารในคลาสนั้นๆ โดยเลือกใช้วลีที่มีค่าความถี่เอกสารมากกว่าค่าที่กำหนดไว้มาใช้

การสร้างจะสร้างทีละวลี ของแต่ละเอกสารในคลาส โดยแทนแต่วลีที่ไม่ซ้ำกันเป็นโหนดของต้นไม้ของคลาส จากนั้น สร้างวลีย่อยจากแต่ละวลี โดยวลีย่อย คือ ลำดับของคำที่ต่อเนื่องกันที่มีขนาดน้อยกว่าของวลี 1 คำ ถ้ายังไม่มีโหนดของวลีย่อยนั้น ก็ทำการเพิ่มโหนดเข้าไปในต้นไม้ของคลาส และสร้างวลีย่อยไปเรื่อยๆ จนสุดท้ายจะได้โหนดใบ คือ คำเดี่ยว และทำทุกวลีที่สกัดได้ของเอกสารจนครบทุกวลี และทุกเอกสารในคลาส

ตัวอย่างเช่น มี 2 วลีที่สกัดได้ คือ “*k-nearest neighbor*” และ “*effici categor algorithm*” นำมาสร้างต้นไม้วลีจะได้ดังรูป 3.3



รูปที่ 3.3 ตัวอย่างต้นไม้วลีที่ได้จากวลี “*k-nearest neighbor*” และ “*effici categor algorithm*”

3.3.2 การลดพีเจอร่วลีระดับ Local

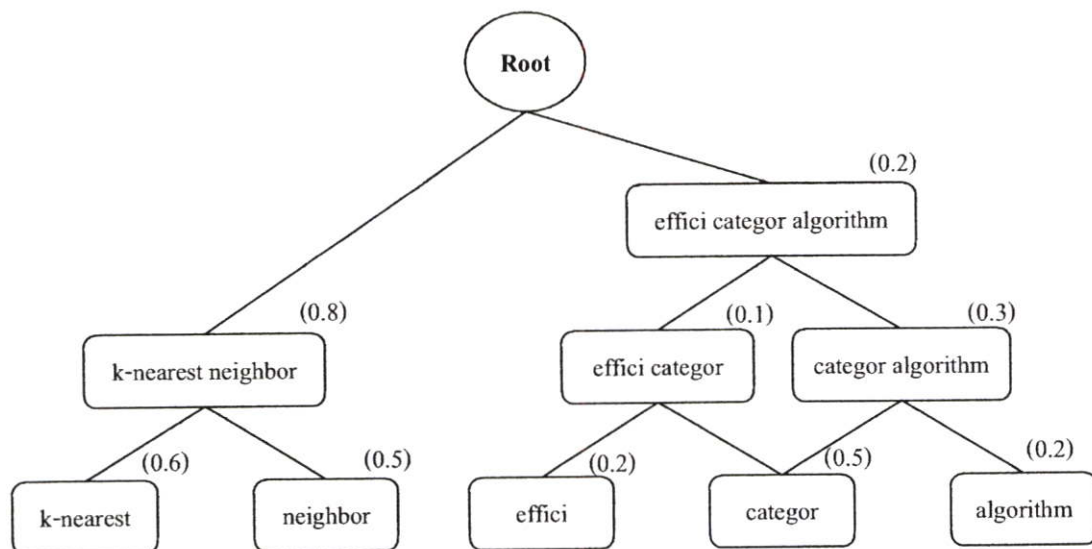
ในระดับ Local นี้ จะเป็นการลดพีเจอร่วลีจากต้นไม้วลีของแต่ละคลาส โดยทำการคำนวณค่าคะแนนของพีเจอร่วลีให้กับแต่ละโหนด ในงานวิจัยนี้ใช้ Odds Ratio Measure [10] โดยค่า Odds Ratio หาได้จากสมการที่ 2.3

นอกจากนี้ยังใช้เกณฑ์จากความถี่เอกสารในการเลือกพีเจอร่วลีอีกด้วย โดยกำหนดค่าต่ำสุดของความถี่เอกสารไว้

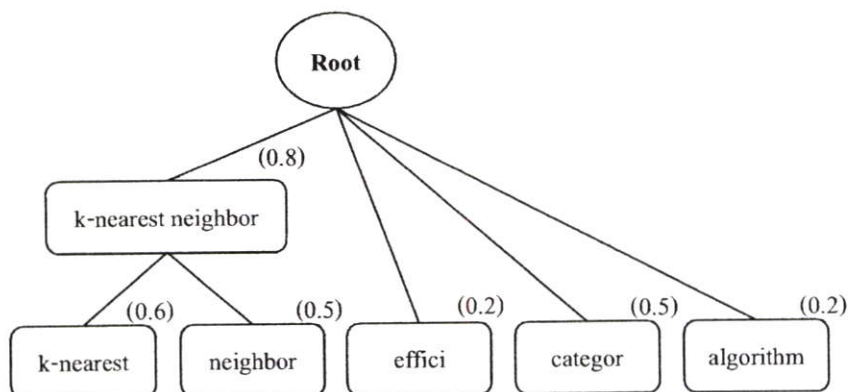
ในการเลือกพีเจอร่วลีนั้นทำได้โดย

1. คำนวณค่า OR ให้กับแต่ละ โหนดของต้นไม้วลีแล้วทำการพิจารณาไปที่ละ โหนด
2. พิจารณาลบโหนดของวลีที่มีค่าคะแนนน้อยกว่าโหนดใดโหนดหนึ่ง ที่เป็นวลีย่อยของมัน ซึ่งรวมถึงโหนดในระดับต่ำสุดที่เป็น โหนดใบด้วย หรือกล่าวได้ว่า ลบโหนดที่มีค่าคะแนนน้อยกว่าค่าคะแนนสูงสุดของวลีที่อยู่ในต้นไม้ย่อย (Sub-tree) ของมัน

จากตัวอย่างต้นไม้วลี รูปที่ 3.3 สมมติว่าทำการคำนวณคะแนนแล้วได้ผลค่า Odds Ratio ดังรูปที่ 3.4 โดยหลังจากผ่านกระบวนการการลดพีเจอร่วลีแล้ว จะได้ต้นไม้วลี ดังแสดงในรูปที่ 3.5



รูปที่ 3.4 ตัวอย่างต้นไม้วลีจากรูปที่ 3.3 ที่ให้คะแนนแล้ว



รูปที่ 3.5 ตัวอย่างต้นไม้วลีจากรูป 3.4 ที่ผ่านการเลือกพีเจอร์แล้ว

3.3.3 การลดพีเจอร์วลีระดับ Global

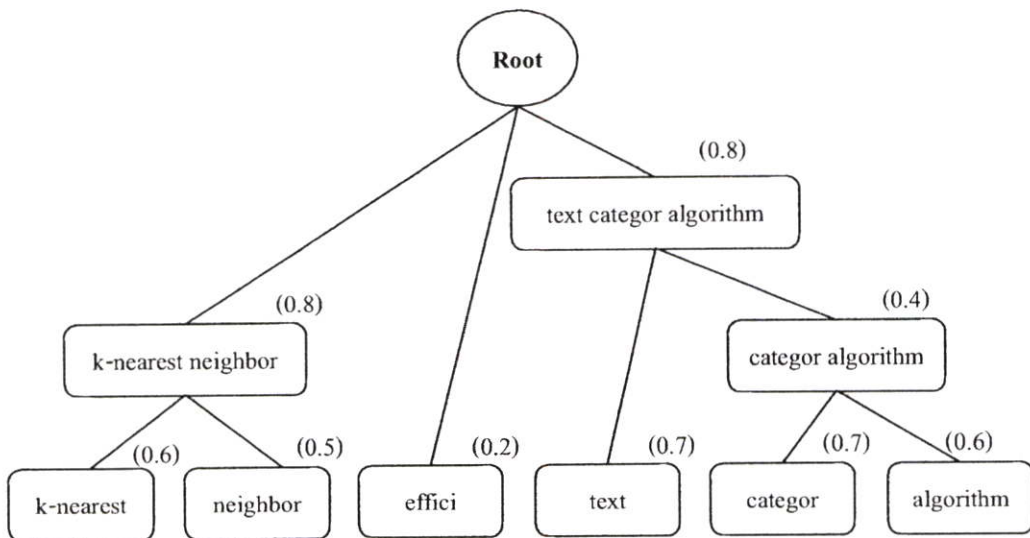
ในการลดพีเจอร์วลีระดับ Global นั้น ต้นไม้วลีของแต่ละคลาสจะถูกนำมารวมกัน โดยในโหนดที่มีร่วมกัน การให้ค่าคะแนนจะเป็น การรวมค่า Odds Ratio ของโหนดนั้นในแต่ละคลาส เพื่อใช้ในการลดพีเจอร์แบบ Global ซึ่งเมื่อได้ต้นไม้วลีจากการรวมต้นไม้ของแต่ละคลาสแล้ว ก็ทำการบวกรวมการลด เช่นเดียวกับในระดับ Local

ตัวอย่างเช่น มีต้นไม้วลีที่ผ่านการเลือกพีเจอร์ระดับ Local ที่ได้จากรูป 3.5 และ c_2 แสดงดังรูป 3.6



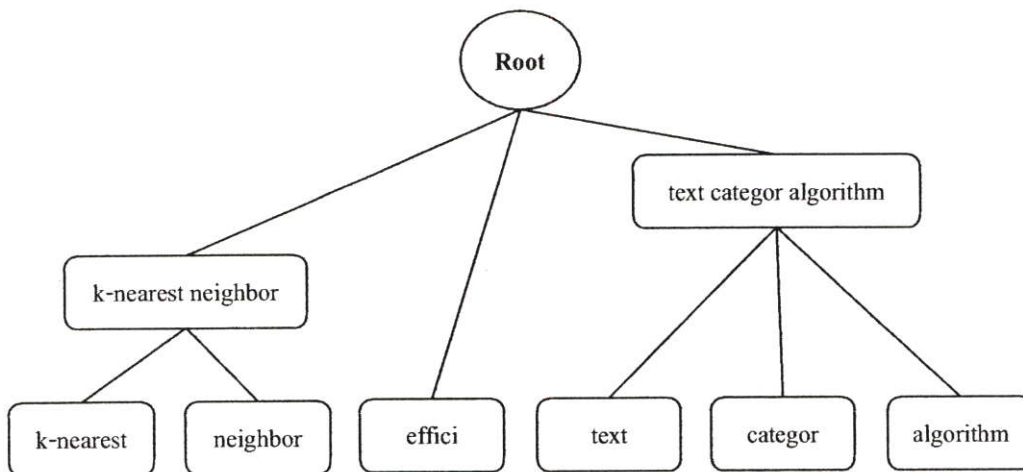
รูปที่ 3.6 ตัวอย่างต้นไม้วลีของคลาส c_2 ที่ผ่านการเลือกพีเจอร์ระดับ Local แล้ว

เมื่อทำการรวมต้นไม้วลีในรูป 3.5 และ 3.6 แล้วได้ต้นไม้วลีของคลาสทั้งหมด โดยโหนดที่มีร่วมกันจะนำค่าคะแนนมาบวกกัน แสดงดังรูปที่ 3.7



รูปที่ 3.7 ต้นไม้วลีจากการรวมกันของต้นไม้ในรูปที่ 3.5 และ 3.6

เมื่อรวมได้เป็นต้นไม้วลีสำหรับการเลือกวลีในระดับ Global แล้ว ทำการลดพีเจอร์อีกครั้ง ด้วยหลักการเดียวกันกับการลดในระดับ Local คือ ลบโหนดที่มีค่าคะแนนน้อยกว่าค่าคะแนนสูงสุดของวลีที่อยู่ในต้นไม้ย่อยของมัน โดยจากรูป 3.7 เมื่อทำการลดพีเจอร์แล้วได้ต้นไม้วลี แล้วได้ต้นไม้วลี ดังรูปที่ 3.8



รูปที่ 3.8 ต้นไม้วลีที่ได้หลังจากทำการลดพีเจอร์ จากต้นไม้วลีในรูปที่ 3.7

สุดท้ายจะได้เซตของพีเจอร์วลี จากวลีในทุกโหนดของต้นไม้วลีที่ได้ ยกเว้นโหนดใบ เพราะเป็นคำเดี่ยว จากตัวอย่าง เราจะได้ เซตของพีเจอร์เป็น

$$Phrase_Features = \{ "k-nearest neighbor", "text categor algorithm" \}$$

3.4 การแทนเอกสาร (Document Representation)

ในการแทนเอกสาร เราจำเป็นต้องแทนเอกสารด้วย คำเดี่ยว และวลี ร่วมกันเพื่อเพิ่มประสิทธิภาพในการจำแนกประเภทเอกสาร ดังนั้น จากต้นไม้วลีสุดท้ายที่ได้ เราจะทำการเลือกพีเจอร์ที่เป็นคำเดี่ยวโดยดูที่เฉพาะ โหนดใบ ซึ่งสามารถใช้เกณฑ์ต่างๆ ในการเลือกได้

ในงานวิจัยของเราเลือกใช้คำเดี่ยวที่มีค่าคะแนนมากกว่าเกณฑ์ที่ตั้งไว้ มาใช้เป็นพีเจอร์ จากตัวอย่างต้นไม้วลีในรูปที่ 3.7 สมมติเราตั้งเกณฑ์คะแนนไว้เป็น 0.3 ดังนั้น เซตของพีเจอร์วลีที่เป็นคำเดี่ยวที่เราได้คือ

$$Word_Features = \{ "k-nearest", "neighbor", "text", "categor", "algorithm" \}$$

สุดท้ายนำพีเจอร์ทั้ง 2 มารวมกัน เพื่อใช้แทนเอกสาร

$$Features = \{ "k-nearest", "neighbor", "text", "categor", "algorithm", "k-nearest neighbor", "text categor algorithm" \}$$

Algorithm: Phrase-tree-based Feature Reduction**Input:** D – training set of documents; $d = d_1, d_2, \dots, d_n$ C – set of categories; $C = c_1, c_2, \dots, c_m$ L – set of label of each documents in D ; $L(d_i) \in C$

min_f – minimum frequency

min_OR – minimum Odds Ratios Score

Output: F : Set of Features $F = \emptyset$ *// Phrase I : Phrase-tree Construction & Local Reduction*For $i=1$ to m $P = \emptyset$ For each document $d_j \in D$ that $L(d_j) = c_i$ PH = set of extracted phrases from d that have document frequency $> \text{min}_f$ $P = P \cup \text{PH}$

End

 $P_tree_i = \text{construct_phrase_tree}(P)$ *//generate phrase-tree rom P* $P_tree_i = \text{calculate_odds_ratio}(P_tree_i)$ *//calculate OR for each node in tree* $P_tree_i = \text{reduce_phrase_tree}(P_tree_i, \text{min_OR})$ *//Local Reduction*

End

// Phrase II & III : Phrase-tree Merging & Global Reduction $M_tree = \emptyset$ For $i=1$ to m *//merge phrase-tree from each category* $M_tree = \text{merge_tree}(M_tree, P_tree_i)$

End

 $M_tree = \text{reduce_phrase_tree}(M_tree, \text{min_OR})$ *//Global Reduction* $F = \text{all nodes from } M_tree$ **รูปที่ 3.9** อัลกอริทึมการลดฟีเจอร์โดยใช้ต้นไม้วลี

ในบทที่ 4 จะกล่าวถึง การทดลองและผลการทดลองกับชุดข้อมูลที่น่ามาทดลองกับ กระบวนการทำงานของงานวิจัยนี้ โดยทดลองเปรียบเทียบประสิทธิภาพในการใช้ฟีเจอร์ในแบบ ต่างๆ

บทที่ 4

การทดลองและผลการทดลอง

ในการพิจารณาประสิทธิภาพของพีเจอร์ เราจะพิจารณาที่ความถูกต้องในการจำแนกประเภทเอกสาร เมื่อแทนเอกสารด้วยพีเจอร์นั้นๆ และพิจารณาที่จำนวนของพีเจอร์ด้วย ซึ่งควรจะมีจำนวนน้อย แต่ให้ความถูกต้องสูง ในการทดลองได้ทำการเปรียบเทียบประสิทธิภาพของพีเจอร์ที่ได้จากต้นไม้วัดกับพีเจอร์แบบอื่นๆ โดยมีรายละเอียดในการทดลองต่อไปนี้

4.1 การวัดความถูกต้องของการจำแนกประเภทเอกสาร

การวัด ความถูกต้องของการจำแนกประเภทเอกสาร ได้ใช้ตัววัดที่ใช้ในงานวิจัยทั่วไป คือ ตัววัดอัตราความถูกต้อง (Accuracy Rate) และตัววัดความถูกต้อง F-measure โดยมีรายละเอียดของแต่ละตัวชี้วัด ดังต่อไปนี้ [13] และ [15]

4.1.1 ตัววัดอัตราความถูกต้อง (Accuracy Rate)

ตัววัดประสิทธิภาพอัตราความถูกต้อง ในการจำแนกประเภทเอกสาร มีรายละเอียดดังสมการที่ (4.1)

$$R = \left(\frac{Ndc}{Ndt} \right) \times 100 \quad (4.1)$$

โดยที่ Ndc คือ จำนวนเอกสารที่ถูกจำแนกประเภทถูกต้อง

Ndt คือ จำนวนเอกสารที่นำมาทดสอบ

4.1.2 ตัววัดความถูกต้อง F-Measure

ตัววัดความถูกต้อง F-measure ใช้ในการวัดความถูกต้องของการจำแนกประเภทเอกสาร โดยจะมีการคำนวณหาค่า Recall และ Precision ค่า F-measure จะมีค่าอยู่ระหว่าง 0 ถึง 1 โดยที่ค่านี้จะบอกว่าผลการจำแนกประเภทเอกสารนั้น มีความถูกต้องมากน้อยเพียงใด ถ้าค่าที่คำนวณได้ มีค่าน้อย แสดงว่ามีความถูกต้องต่ำ ถ้าค่าที่คำนวณได้มีค่าเป็น 1 แสดงว่ามีความถูกต้องสูงสุด การคำนวณหาค่า F-measure $F(i)$ ของประเภทเอกสาร ลำดับที่ i จะต้องคำนวณหาค่า Recall และ Precision ซึ่งส่วนที่นำมาหาค่า Recall และ Precision พิจารณาจากรูปที่ 4.1

สมการที่ใช้สำหรับหาค่า Recall และ Precision คือสมการที่ (4.2) และ (4.3) ตามลำดับ

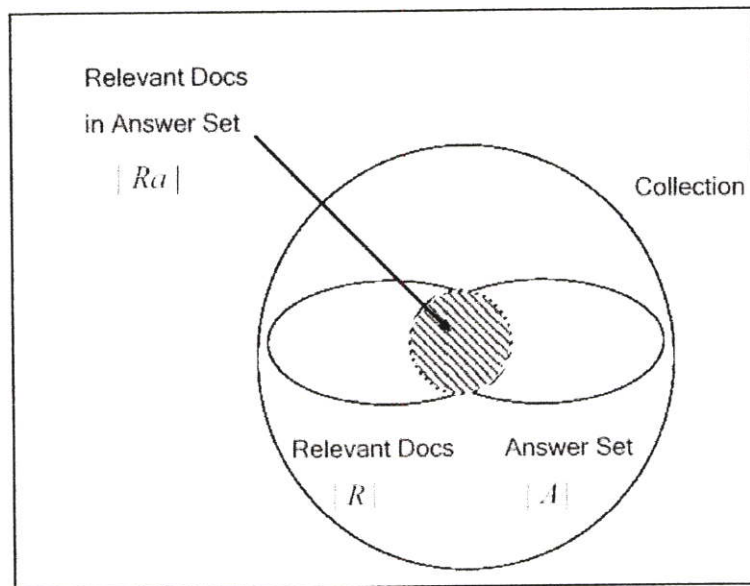
$$Recall(i) = \frac{|Ra_i|}{|R_i|} \quad (4.2)$$

$$Precision(i) = \frac{|Ra_i|}{|A_i|} \quad (4.3)$$

โดยที่ $|Ra_i|$ คือ จำนวนของเอกสารในชุดทดสอบที่อยู่ในประเภทที่ i และเมื่อทำการจำแนกประเภทแล้ว ก็ได้ผลการจำแนกประเภทว่าอยู่ในประเภทที่ i

$|R_i|$ คือ จำนวนของเอกสารในชุดทดสอบที่อยู่ในประเภทที่ i

$|A_i|$ คือ จำนวนของเอกสารที่ได้จากการจำแนกประเภทว่าอยู่ในประเภทที่ i



รูปที่ 4.1 แสดงส่วนที่ใช้วัดค่า Precision และ Recall ในการร้องขอสารสนเทศ [15]

การคำนวณหาค่า $F(i)$ ของ ประเภทเอกสารลำดับที่ i สามารถคำนวณได้จากสมการที่

(4.4)

$$F(i) = \frac{2 \times Recall(i) \times Precision(i)}{Recall(i) + Precision(i)} \quad (4.4)$$

การคำนวณค่า F-measure ของทุกๆ ประเภทเอกสาร สามารถคำนวณได้จากสมการที่ (4.5)

$$F = \sum_{i=1}^{Nc} \frac{|R_i|}{Ndt} F(i) \quad (4.5)$$

โดยที่ Nc คือ จำนวนของประเภท

Ndt คือ จำนวนเอกสารที่นำมาทดสอบ

4.2 การออกแบบการทดลอง

4.2.1 การสุ่มตัวอย่างเอกสาร (Document Sampling)

เนื่องจากในบางชุดเอกสารมีเอกสารจำนวนมาก ทำให้มีความซับซ้อนในการคำนวณสูง ดังนั้นจึงจำเป็นต้องลดจำนวนเอกสารลง โดยนำเทคนิคการสุ่มตัวอย่างในทางสถิติมาใช้ในการสุ่มเอกสารตัวอย่างขึ้นมา เพื่อใช้ในการเรียนรู้และทดสอบต่อไป

รศ.สุรินทร์ นิยมางกูร กล่าวในหนังสือเทคนิคการสุ่มตัวอย่าง [20] ว่า ตัวอย่าง (Sample) ที่ดีคือ ตัวอย่างที่สามารถให้รายละเอียดข้อเท็จจริงเกี่ยวกับคุณลักษณะของประชากรได้เป็นอย่างดี ตามหลักวิชาการทางสถิติมีความเชื่อว่า ตัวอย่างที่เลือกมาโดยวิธีการสุ่มจะไม่ทำให้เกิดความเอนเอียงในการเก็บรวบรวมรายละเอียดข้อเท็จจริงเกี่ยวกับคุณลักษณะของประชากร ดังนั้นตัวอย่างที่ดีจึงควรเป็นตัวอย่างที่เลือกมาโดยวิธีการสุ่ม หรือ อาจจะกล่าวได้ว่าตัวอย่างที่ดีก็คือ ตัวอย่างสุ่ม นั่นเอง

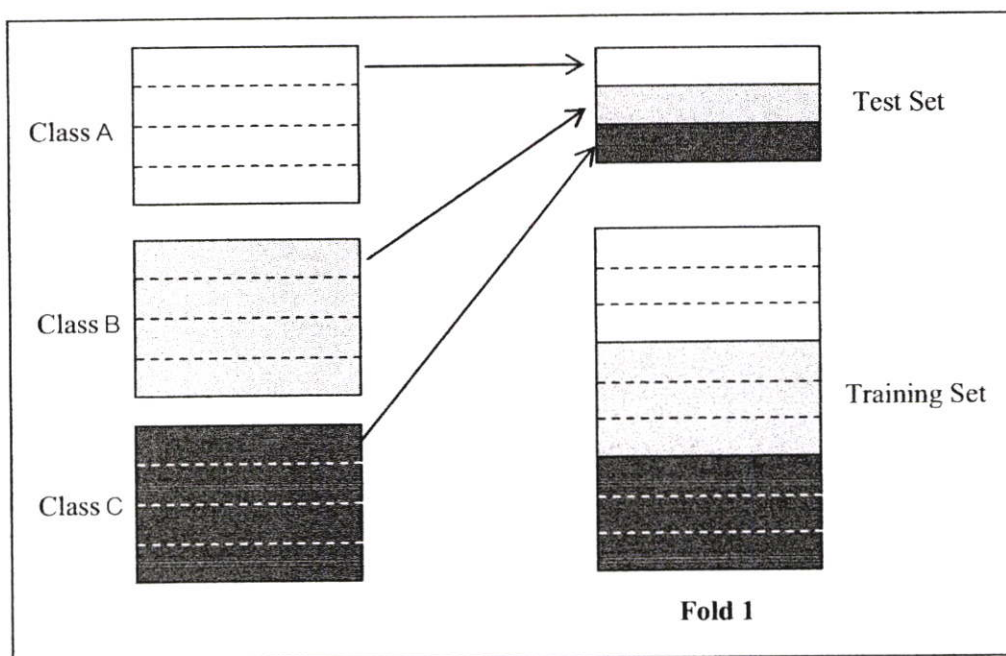
การสุ่มตัวอย่างแบ่งเป็น 2 แบบใหญ่ๆ คือ การสุ่มตัวอย่างที่ไม่เป็นไปตามโอกาสทางสถิติ (Non-Probability Sampling) และ การสุ่มตัวอย่างที่เป็นไปตามโอกาสทางสถิติ (Probability Sampling) โดยการสุ่มตัวอย่างที่คำนึงถึงความน่าจะเป็นในการสุ่ม มี 4 วิธี คือ การสุ่มตัวอย่างแบบง่าย (Simple Random Sampling) การสุ่มอย่างมีระบบ (Systematic Random Sampling) การสุ่มตัวอย่างแบบเป็นชั้นภูมิ (Stratified Random Sampling) และ การสุ่มแบบแบ่งกลุ่ม (Cluster Sampling)

การสุ่มตัวอย่างแบบเป็นชั้นภูมิ มีการจัดแบ่งประชากรเป็นกลุ่มหรือชั้นย่อยๆ ก่อน แล้วเลือกสุ่มตัวอย่างตามสัดส่วน (Proportional) ในแต่ละชั้น จากนั้นจึงใช้การสุ่มอย่างง่าย เช่น แบ่งนักศึกษาตามคณะต่างๆ หาขนาดกลุ่มตัวอย่าง จากนั้นเทียบสัดส่วนตามขนาด แล้วจับฉลาก เป็นต้น วัตถุประสงค์หลักของการสุ่มตัวอย่างแบบแบ่งเป็นชั้นภูมิคือ เพื่อให้ได้กลุ่มตัวอย่างที่มีองค์ประกอบของลักษณะต่าง ๆ ใกล้เคียงกับกลุ่มประชากร และให้ได้กลุ่มตัวอย่างที่สามารถตอบสนองวัตถุประสงค์ของการวิจัยได้

ในงานวิจัยนี้เลือกประยุกต์ใช้ การสุ่มตัวอย่างแบบเป็นชั้นภูมิ เพื่อรักษาลักษณะของคำภายในแต่ละประเภทเอกสารไว้ โดยในขั้นแรกจะแบ่งชุดเอกสารตามประเภทเอกสาร และในแต่ละประเภทเอกสารทำการแบ่งชั้นภูมิตามการกระจายตัวของเอกสารในคำต่างๆ ชั้นภูมิในแต่ละประเภทเอกสารของเราในที่นี้ก็คือ คำแต่ละคำ ซึ่งมีจำนวนชั้นภูมิเท่ากับจำนวนคำทั้งหมด ของประเภทเอกสารนั้น เมื่อแบ่งชั้นภูมิแล้วก็ทำการสุ่มเอกสาร จากแต่ละชั้นภูมิ โดยใช้การสุ่มอย่างง่าย จำนวนเอกสารที่จะสุ่มขึ้นมานั้น เป็นไปตามสัดส่วนที่เรากำหนดไว้ เช่น กำหนดไว้ที่ 0.1 ถ้าในคำที่ 1 มีจำนวนเอกสารเท่ากับ 100 จำนวนตัวอย่างที่เราจะสุ่มมาในชั้นภูมิคำนี้จะเท่ากับ 1 เอกสาร

4.2.2 การตรวจสอบไขว้ (Cross-Validation)

การตรวจสอบไขว้ เป็นวิธีการทดสอบการจำแนกประเภทเอกสารที่นักวิจัยนิยมใช้ โดย k-fold Cross Validation หมายถึง ทำการแบ่งชุดเอกสารออกเป็น k ชุด แล้วทำการทดสอบ โดยแต่ละครั้ง ใช้ 1 ส่วนเป็นชุดทดสอบส่วน k-1 ชุดที่เหลือเป็นชุดทดสอบ สับเปลี่ยนชุดทดสอบจนใช้ทุกส่วนเป็นชุดทดสอบ ซึ่งเท่ากับทำการทดลอง k ครั้ง แล้วนำผลที่วัดได้จากทั้ง k ชุดการทดลองมาหาค่าเฉลี่ย ตัวอย่างเช่น การทำ 4-fold Cross validation เริ่มโดยทำการสุ่มแบ่งเอกสารในแต่ละประเภทเอกสารออกเป็น 4 ส่วน แล้วนำ 1 ส่วนจากแต่ละประเภทมารวมเป็นกลุ่มสำหรับทดสอบ แล้วใช้ 3 ส่วนที่เหลือจากแต่ละประเภทมารวมกันเป็นกลุ่มเรียนรู้ แสดงดังรูปที่ 4.2



รูปที่ 4.2 แสดงตัวอย่างการสร้างชุดเอกสารเรียนรู้และทดสอบจากกลุ่มเอกสาร 3 ประเภท A B และ C สำหรับการทำ 4-fold cross validation ในชุดการทดลองลำดับที่ 1 (fold 1)

4.2.3 รูปแบบของพีเจอร์ที่ใช้ในการทดลอง

ในการทดลองจำแนกประเภทเอกสารได้ใช้พีเจอร์แบบต่างๆ ดังนี้

- คำเดี่ยว
- คำเดี่ยวและไบนารี-แกรม โดยสร้างไบนารี-แกรมจากลำดับของคำในแต่ละวลีที่สกัดได้
- คำเดี่ยวและวลี โดยวลีที่ใช้คือวลีที่สกัดได้
- พีเจอร์จากต้นไม้วลี คือวิธีการที่นำเสนอในงานวิจัยนี้ โดยพีเจอร์จะถูกสร้างและทำการเลือกโดยใช้ต้นไม้วลี

4.3 ชุดเอกสารที่ใช้ในการทดลอง

เอกสารที่ใช้ในการทดลองมีจำนวน 2 ชุดเอกสารใหญ่ คือ

ก) ชุดเอกสาร PDDP [16] เป็นเอกสารเว็บ ซึ่งแบ่งเป็น 3 ชุดย่อย F-series, J-series และ K-series โดยแต่ละเอกสารอยู่ในรูปแบบเอกสาร HTML เป็นชุดเอกสารที่ถูกนำมาใช้ในงานวิจัยทางการจำแนกประเภทเอกสาร

ข) ชุดข้อมูล Reuters-Top10 กัดเลือกมาจาก ชุดข้อมูลข่าว Reuters-21578 [18] 10 กลุ่มแรกที่มีจำนวนข่าวมากที่สุดจาก 135 กลุ่มเอกสารข่าว โดยแต่ละเอกสารอยู่ในรูปแบบ XML ซึ่งชุดข้อมูล Reuters-Top10 [19] เป็นชุดข้อมูลที่ใช้ในการทดสอบการจำแนกประเภทเอกสารและการจัดกลุ่มเอกสาร ที่ใช้กันอย่างแพร่หลาย

ชุดเอกสารทั้ง หกจะผ่านกระบวนการเตรียมเอกสารก่อน ตัวอย่างเอกสาร และการเตรียมเอกสารของเอกสารแต่ละชุด แสดงในภาคผนวก ก

4.4 ชุดเอกสาร PDDP F-series ที่ใช้ในการทดลองและผลการทดลอง

4.4.1 ชุดเอกสาร PDDP F-series ที่ใช้ในการทดลอง

ชุดเอกสาร PDDP F-series มีเอกสารทั้งหมด 98 แบ่งเป็น 4 ประเภท อยู่ในรูปแบบแฟ้มข้อมูล HTML ซึ่งมีรายละเอียดอยู่ในภาคผนวก ก รายละเอียดของแต่ละ ประเภทเอกสารแสดงในตารางที่ 4.1 ในการทดลองจำแนกประเภทเอกสาร ได้แบ่งชุดเอกสาร PDDP F-series เพื่อใช้ในการเรียนรู้ และใช้ในการทดสอบ สำหรับทำ 4-fold cross validation ได้จำนวนเอกสารสำหรับเรียนรู้และทดสอบในแต่ละกลุ่มการทดลองดังแสดงในตาราง 4.2

4.4.2 ผลการทดลองของชุด เอกสาร PDDP F-series เมื่อใช้พีเจอร์ที่ได้จากการเลือกพีเจอร์โดยใช้ต้นไม้วลี

สร้างต้นไม้วลีและทำการเลือกพีเจอร์โดยใช้ต้นไม้วลีโดยกำหนดให้แต่ละเทอมมีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 1 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.1 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอร์ดังกล่าว ด้วยวิธี เค-เนียร์สเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.2 และ 4.3 ตามลำดับ รายละเอียดของพีเจอร์ที่ได้แสดงในภาคผนวก ข

ตารางที่ 4.1 แสดงจำนวนเอกสาร F-series แบ่งตามประเภทเอกสาร

ชื่อประเภท	จำนวนเอกสาร
Business & Finance (B)	22
Electronic Communication (E)	23
Labor (L)	26
Manufacturing (M)	27
รวม	98

ตารางที่ 4.2 แสดงจำนวนเอกสาร F-series ที่ใช้ในการเรียนรู้และทดสอบในแต่ละชุดการทดลอง

ชื่อประเภท (ตัวย่อ)	จำนวนเอกสารในชุดการทดลอง (fold) ที่							
	1		2		3		4	
	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ
B	16	6	16	6	17	5	17	5
E	17	6	17	6	17	6	18	5
L	20	6	20	6	19	7	19	7
M	21	6	20	7	20	7	20	7
รวม	74	24	73	25	73	25	74	24
	98		98		98		98	

ตารางที่ 4.3 แสดงจำนวนฟีเจอร์จากต้นไม้วลี ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลองของชุดเอกสาร PDDP F-series

ชุดที่	จำนวนฟีเจอร์จากต้นไม้วลี	
	ทั้งหมด	ฟีเจอร์ที่ถูกเลือก
1	23,469	2,832
2	23,861	2,797
3	22,911	2,738
4	23,415	2,740
ค่าเฉลี่ย	23,414	2,777

ตารางที่ 4.4 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียบเรสเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์จากต้นไม้วลี ของชุดเอกสาร PDDP F-series

ชุด ที่	k									ค่า สูงสุด
	1	2	3	4	5	6	7	8	9	
1	100.0	100.0	95.83	100.0	91.67	95.83	87.50	95.83	95.83	100.0
2	97.10	88.00	97.10	96.00	96.00	96.00	96.00	92.00	96.00	97.10
3	96.00	96.00	100.0	100.0	100.0	100.0	96.00	100.0	96.00	100.0
4	92.88	91.67	92.88	87.50	87.50	87.50	87.50	87.50	83.33	92.88
ค่าเฉลี่ย	96.50	93.92	96.45	95.88	93.79	94.83	91.75	93.83	92.79	97.50

ตารางที่ 4.5 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เนียบเรสเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์จากต้นไม้วลี ของชุดเอกสาร F-series

ชุด ที่	k									ค่า สูงสุด
	1	2	3	4	5	6	7	8	9	
1	1.0000	1.0000	0.9580	1.0000	0.9161	0.9580	0.8723	0.9580	0.9580	1.0000
2	0.9710	0.8709	0.9710	0.9600	0.9597	0.9597	0.9597	0.9224	0.9597	0.9710
3	0.9593	0.9600	1.0000	1.0000	1.0000	1.0000	0.9595	1.0000	0.9595	1.0000
4	0.9176	0.9095	0.9176	0.8713	0.8713	0.8713	0.8713	0.8713	0.8171	0.9176
ค่าเฉลี่ย	0.9620	0.9351	0.9617	0.9578	0.9368	0.9473	0.9157	0.9379	0.9236	0.9722

4.4.3 ผลการทดลองของชุดเอกสาร PDDP F-series เมื่อใช้ฟีเจอร์แบบคำเดียว

สร้างฟีเจอร์แบบคำเดียว จากทุกคำในเอกสารในชุดเรียนรู้ แล้วทำการเลือกฟีเจอร์ที่มีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 1 เอกสาร และทำการทดลองจำแนกประเภทเอกสารโดยในแต่ละชุดการทดลองได้จำนวนฟีเจอร์จากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.6 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยฟีเจอร์ดังกล่าว ด้วยวิธี เค-เนียบเรสเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.7 และ 4.8 ตามลำดับ

ตารางที่ 4.6 แสดงจำนวนฟีเจอร์แบบคำเดี่ยวที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง
ของชุดเอกสาร PDDP F-series

ชุดที่	จำนวนฟีเจอร์แบบคำเดี่ยว	
	ทั้งหมด	ฟีเจอร์ที่ถูกเลือก
1	5,861	2,676
2	5,948	2,664
3	5,680	2,597
4	5,943	2,678
ค่าเฉลี่ย	5,858	2,654

ตารางที่ 4.7 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียร์เซนเบอร์ ค่า $k = 1$
ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์แบบคำเดี่ยว ของชุดเอกสาร PDDP F-series

ชุด ที่	k									ค่า สูงสุด
	1	2	3	4	5	6	7	8	9	
1	100.0	100.0	91.67	95.83	91.67	91.67	87.50	87.50	87.50	100.0
2	96.00	96.00	96.00	96.00	92.00	92.00	96.00	96.00	96.00	96.00
3	100.0	100.0	100.0	100.0	100.0	100.0	96.00	96.00	96.00	100.0
4	91.67	87.50	91.67	87.50	91.67	87.50	87.50	87.50	83.33	91.67
ค่าเฉลี่ย	96.92	95.88	94.83	94.83	93.83	92.79	91.75	91.75	90.71	96.92

ตารางที่ 4.8 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เนียร์เซนเบอร์
ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์แบบคำเดี่ยว ของชุดเอกสาร F-
series

ชุด ที่	k									ค่า สูงสุด
	1	2	3	4	5	6	7	8	9	
1	1.0000	1.0000	0.9161	0.9580	0.9161	0.9161	0.8723	0.8723	0.8723	1.0000
2	0.9600	0.9600	0.9600	0.9600	0.9224	0.9224	0.9597	0.9597	0.9597	0.9600
3	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.9595	0.9595	0.9595	1.0000
4	0.9095	0.8713	0.9095	0.8713	0.9095	0.8713	0.8713	0.8713	0.8171	0.9095
ค่าเฉลี่ย	0.9674	0.9578	0.9464	0.9473	0.9370	0.9274	0.9157	0.9157	0.9022	0.9674

4.4.4 ผลการทดลองของชุดเอกสาร PDDP F-series เมื่อใช้พีเจอร์แบบคำเดียว และไป-แกรม

สร้างพีเจอร์ที่เป็นคำเดียว และไป-แกรม โดยสร้างไป-แกรมจากวลีที่สกัดได้ แล้วทำการเลือกพีเจอร์ที่มีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 1 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.9 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอร์ดังกล่าว ด้วยวิธี เค-เนียร์สเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.10 และ 4.11 ตามลำดับ

ตารางที่ 4.9 แสดงจำนวนคำเดียวและไป-แกรมที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลองของชุดเอกสาร PDDP F-series

ชุดที่	จำนวนพีเจอร์คำเดียวและไป-แกรม	
	ทั้งหมด	พีเจอร์ที่ถูกเลือก
1	17,750	3,473
2	17,944	3,444
3	17,259	3,379
4	17,762	3,413
ค่าเฉลี่ย	17,679	3,427

ตารางที่ 4.10 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า k = 1 ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์แบบคำเดียวและไป-แกรม ของชุดเอกสาร PDDP F-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	100.0	100.0	91.67	100.0	91.67	95.83	87.50	91.67	87.50	100.0
2	96.00	92.00	96.00	96.00	92.00	92.00	96.00	96.00	96.00	96.00
3	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	96.00	100.0
4	91.67	87.50	91.67	87.50	87.50	87.50	87.50	87.50	83.33	91.67
ค่าเฉลี่ย	96.92	94.88	94.83	95.88	92.79	93.83	92.75	93.79	90.71	96.92

ตารางที่ 4.11 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เน็ยเรสเนเบอร์
ค่า $k = 1$ ถึง 9 ในแต่ละชุดการเมื่อใช้พีเจอรแบบคำเดี่ยวและ ไบ-แกรม ของชุดเอกสาร
F-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	1.0000	1.0000	0.9161	1.0000	0.9161	0.9580	0.8723	0.9143	0.8723	1.0000
2	0.9600	0.9164	0.9600	0.9600	0.9224	0.9224	0.9597	0.9597	0.9597	0.9600
3	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.9595	1.0000
4	0.9095	0.8713	0.9095	0.8713	0.8713	0.8713	0.8713	0.8713	0.8171	0.9095
ค่าเฉลี่ย	0.9674	0.9469	0.9464	0.9578	0.9274	0.9379	0.9258	0.9363	0.9022	0.9674

4.4.5 ผลการทดลองของชุดเอกสาร PDDP F-series เมื่อใช้พีเจอรคำเดี่ยวและวลี

สร้างพีเจอรจากวลีที่สกัดได้และคำเดี่ยว แล้วทำการเลือกพีเจอรที่มีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 1 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอรจากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.12 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอรดังกล่าว ด้วยวิธีเค-เน็ยเรสเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.13 และ 4.14 ตามลำดับ

4.4.6 เปรียบเทียบผลการทดลองเมื่อใช้พีเจอรแบบต่าง ของชุดเอกสาร F-series

รูปที่ 4.3 ถึง 4.15 แสดงการเปรียบเทียบผลการทดลองเมื่อใช้พีเจอรแต่ละแบบ โดยเปรียบเทียบทั้งในด้านจำนวนพีเจอร อัตราความถูกต้อง และค่า F-measure

ตารางที่ 4.12 แสดงจำนวนพีเจอรคำเดี่ยวและวลี ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง
ของชุดเอกสาร PDDP F-series

ชุดที่	จำนวนพีเจอรคำเดี่ยวและวลี	
	ทั้งหมด	พีเจอรที่ถูกเลือก
1	15,569	3,270
2	15,618	3,222
3	15,159	3,185
4	15,576	3,220
ค่าเฉลี่ย	15,481	3,224

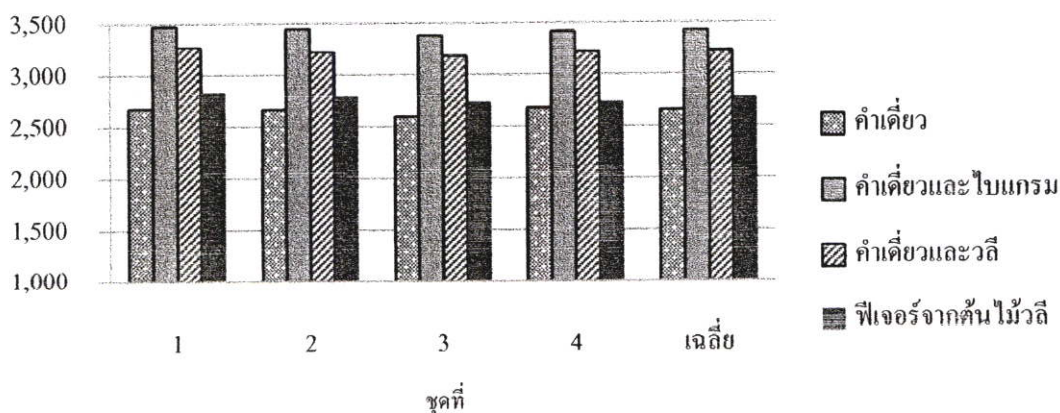
ตารางที่ 4.13 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์ท่าเดียวและวลี ของชุดเอกสาร PDDP F-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	100.0	100.0	95.83	95.83	95.83	95.83	95.83	95.83	95.83	100.0
2	96.00	92.00	96.00	96.00	96.00	96.00	96.00	96.00	96.00	96.00
3	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
4	87.50	87.50	87.50	87.50	87.50	87.50	87.50	87.50	87.50	87.50
ค่าเฉลี่ย	95.88	94.88	94.83	94.83	94.83	94.83	94.83	94.83	94.83	95.88

ตารางที่ 4.14 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์ท่าเดียวและวลี ของชุดเอกสาร F-series

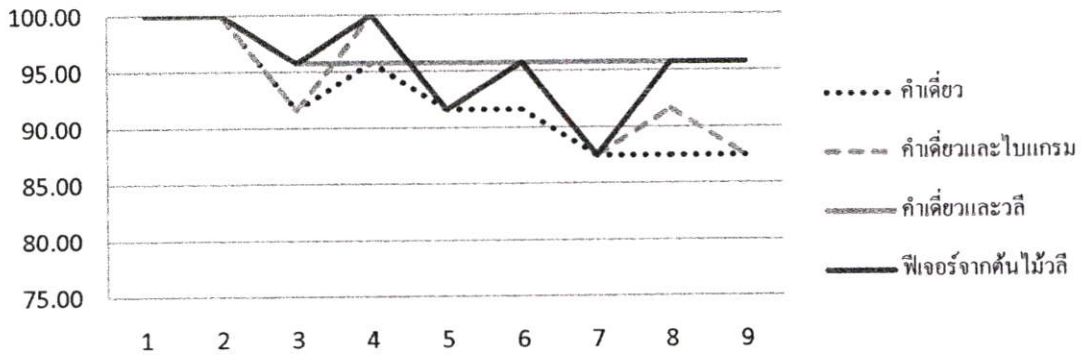
ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	1.0000	1.0000	0.9580	0.9580	0.9580	0.9580	0.9580	0.9580	0.9580	1.0000
2	0.9600	0.9164	0.9600	0.9600	0.9600	0.9600	0.9600	0.9600	0.9224	0.9600
3	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000	0.9595	1.0000
4	0.8713	0.8713	0.8713	0.8713	0.8713	0.8713	0.8713	0.8713	0.8713	0.8713
ค่าเฉลี่ย	0.9578	0.9469	0.9473	0.9473	0.9473	0.9473	0.9473	0.9473	0.9278	0.9578

จำนวนพีเจอร์ท่าที่ใช้ของชุดเอกสาร F-series



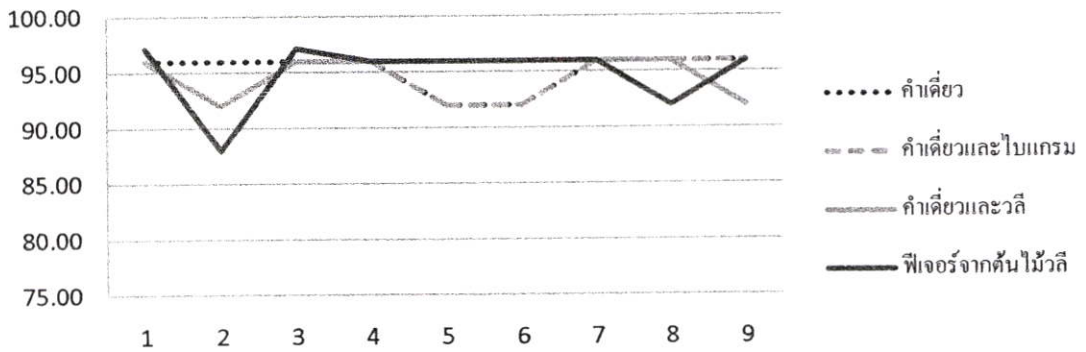
รูปที่ 4.3 แผนภูมิแท่งเปรียบเทียบจำนวนพีเจอร์ท่าที่ใช้ในแต่ละแบบของชุดเอกสาร F-series

อัตราความถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร F-series ชุดที่ 1



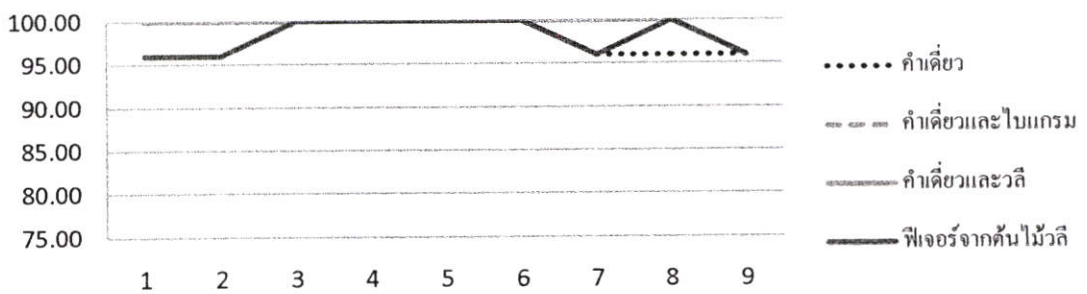
รูปที่ 4.4 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 1

อัตราความถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร F-series ชุดที่ 2



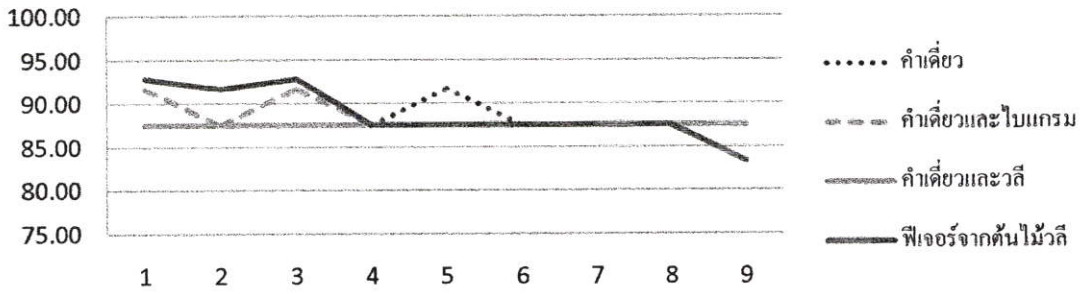
รูปที่ 4.5 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 2

อัตราความถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร F-series ชุดที่ 3



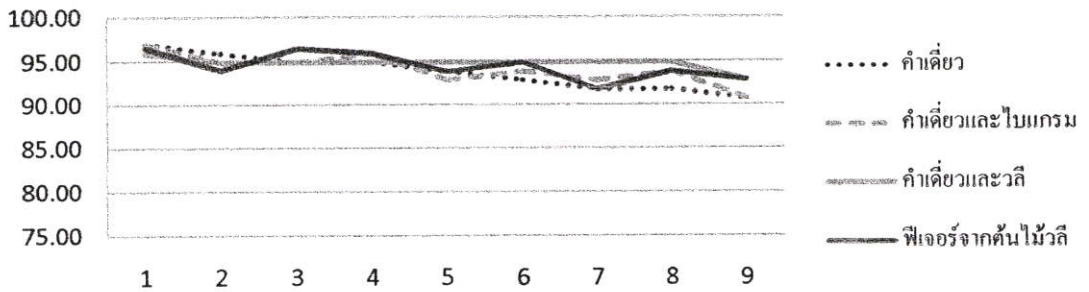
รูปที่ 4.6 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 3

อัตราความถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร F-series ชุดที่ 4



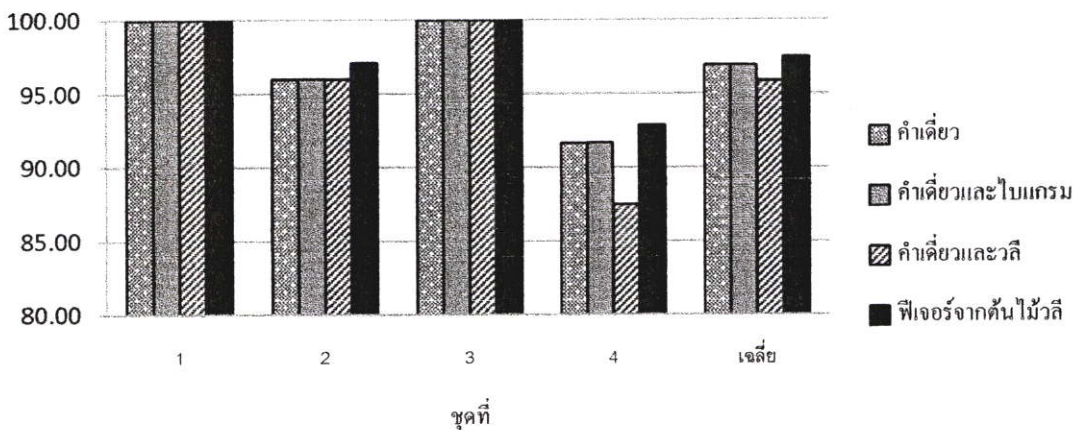
รูปที่ 4.7 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 4

อัตราความถูกต้องเฉลี่ยเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร F-series

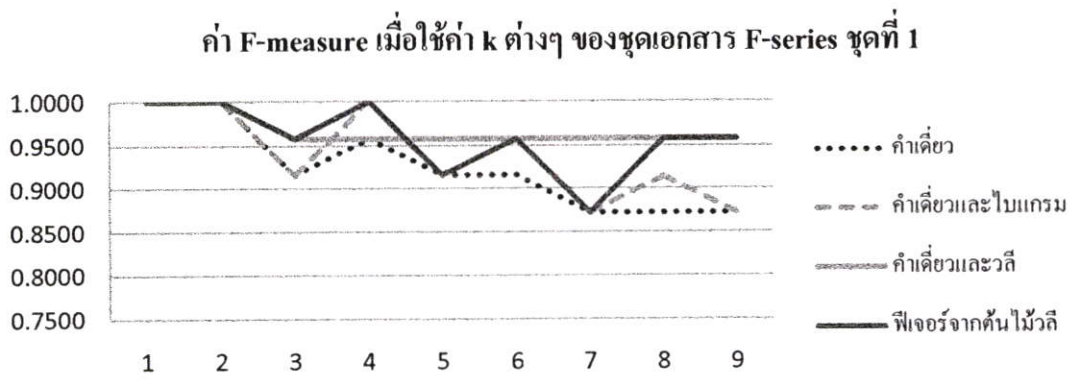


รูปที่ 4.8 กราฟเปรียบเทียบอัตราความถูกต้องเฉลี่ยของทุกชุดการทดลองในแต่ละแบบของชุดเอกสาร F-series

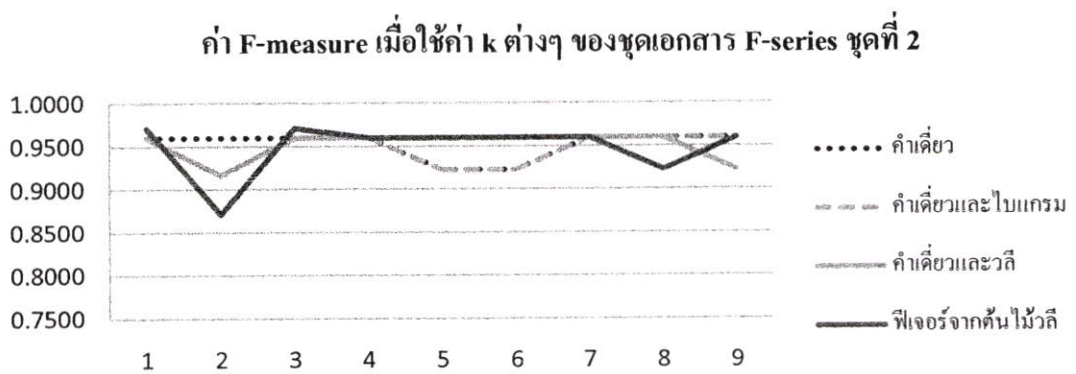
อัตราความถูกต้องสูงสุดของชุดเอกสาร F-series



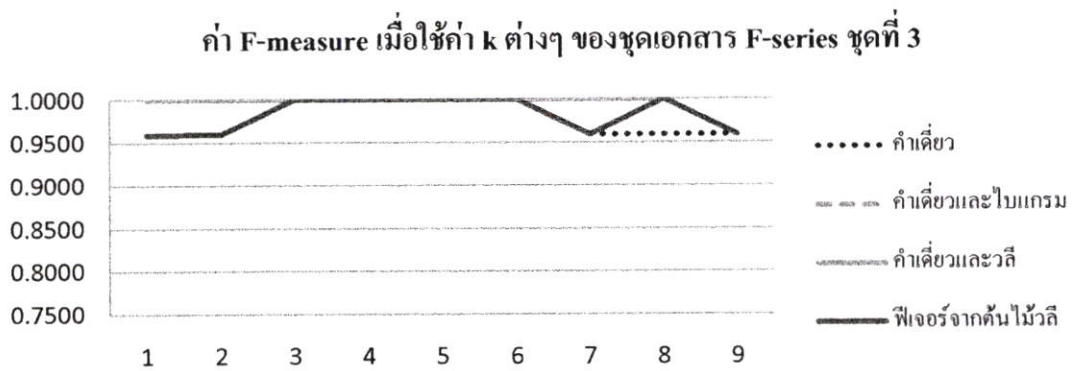
รูปที่ 4.9 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุดในแต่ละแบบของชุดเอกสาร F-series



รูปที่ 4.10 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 1

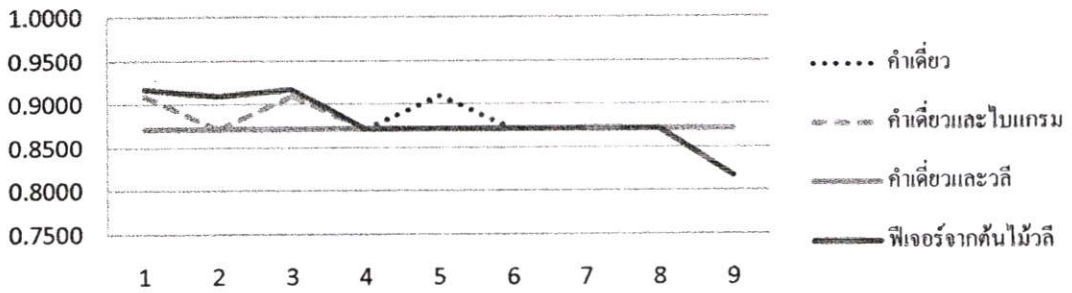


รูปที่ 4.11 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 2



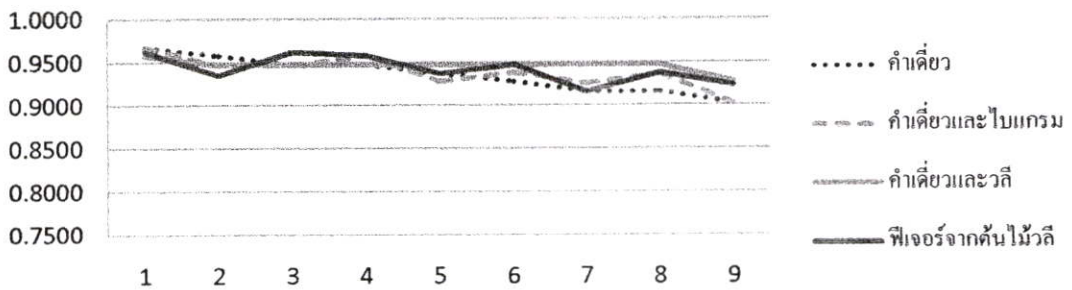
รูปที่ 4.12 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 3

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร F-series ชุดที่ 4



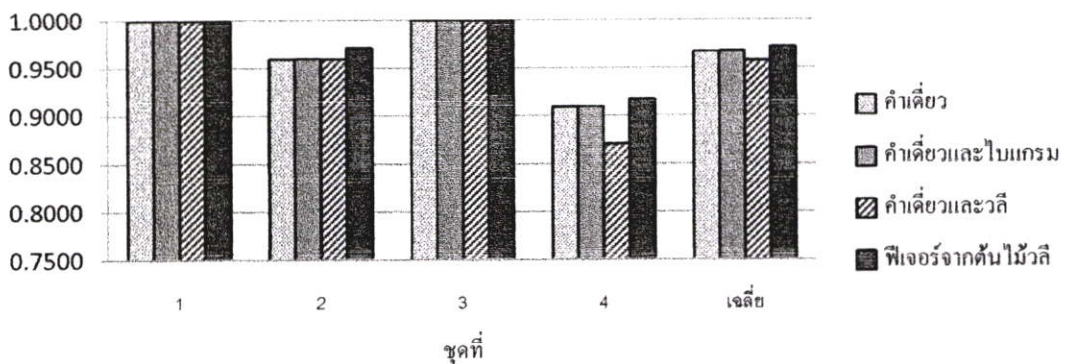
รูปที่ 4.13 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร F-series ชุดที่ 4

ค่า F-measure เฉลี่ยเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร F-series



รูปที่ 4.14 กราฟเปรียบเทียบ ค่า F-measure เฉลี่ยในแต่ละแบบของชุดเอกสาร F-series

ค่า F-measure สูงสุดของชุดเอกสาร F-series



รูปที่ 4.15 แผนภูมิแท่งเปรียบเทียบ ค่า F-measure สูงสุดในแต่ละแบบของชุดเอกสาร F-series

จากรูปที่ 4.3 จะเห็นได้ว่าฟีเจอร์ที่ได้จากต้นไม้วลีมีจำนวนน้อยกว่าฟีเจอร์แบบอื่นๆ ยกเว้นแบบคำเดียว และเมื่อพิจารณาถึงความถูกต้องในการจำแนกประเภทเอกสาร จากรูปที่ 4.9 และ 4.15

พีเจอร်จากต้นไม้วลีให้อัตราความถูกต้องและค่า F-measure เมื่อเฉลี่ยจากทั้ง 4 ชุดสูงกว่าพีเจอร်แบบอื่นๆ ซึ่งเป็นไปตามเป้าหมายของเราที่ต้องการได้พีเจอร်ที่มีประสิทธิภาพ โดยใช้วลีร่วมกับคำเดี่ยวที่ให้ความถูกต้องในการจำแนกประเภทเอกสารสูงแต่มีจำนวนน้อย

4.5 ชุดเอกสาร PDDP J-series ที่ใช้ในการทดลองและผลการทดลอง

4.5.1 ชุดเอกสาร PDDP J-series ที่ใช้ในการทดลอง

ชุดเอกสาร PDDP J-series มีเอกสารทั้งหมด 185 แบ่งเป็น 10 ประเภท ซึ่งในเอกสารชุดนี้ต่างจาก J-series ตรงที่ประเภทของเอกสารจะถูกแบ่งย่อยมากกว่า รายละเอียดของแต่ละประเภทเอกสารแสดงในตารางที่ 4.15

ตารางที่ 4.15 แสดงจำนวนเอกสาร J-series แบ่งตามประเภท

ชื่อประเภท	จำนวนเอกสาร	ชื่อประเภท	จำนวนเอกสาร
Affirmative Action (A)	20	Employee Rights (L)	16
Business Capital (B)	19	Materials Processing (M)	17
Information Systems (C)	19	Personal Management (P)	19
Electronic Commerce (E)	19	Manufacturing Systems (S)	18
Intellectual Property (I)	19	Industrial Partnership (Z)	19

ในการทดลองจำแนกประเภทเอกสารได้แบ่งชุดเอกสาร PDDP J-series เพื่อใช้ในการเรียนรู้ และใช้ในการทดสอบ สำหรับทำ 4-fold cross validation โดยการสุ่มเลือกเอกสารได้จำนวนเอกสารสำหรับเรียนรู้และทดสอบในแต่ละชุดการทดลองดังแสดงในตาราง 4.16 และทดลองเพื่อวัดประสิทธิภาพของพีเจอร်ที่ได้จากการเลือกโดยใช้ต้นไม้วลีในการจำแนกประเภทเอกสารเปรียบเทียบกับพีเจอร်แบบอื่นๆ โดยทำการทดลองสร้างพีเจอร် 4 แบบ สำหรับการจำแนกประเภทคือ 1.พีเจอร်จากต้นไม้วลี 2.คำเดี่ยว 3.คำเดี่ยวและไบน-แกรม และ 4.คำเดี่ยวและวลี

4.5.2 ผลการทดลองของชุดเอกสาร PDDP J-series เมื่อใช้พีเจอร်ที่ได้จากการเลือกพีเจอร်โดยใช้ต้นไม้วลี

สร้างต้นไม้วลีและทำการเลือกพีเจอร်โดยใช้ต้นไม้วลีโดยกำหนดให้แต่ละเทอมมีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 2 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอร်จากเอกสารชุดเรียนรู้ดังแสดงในตารางที่ 4.17 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วย

พีเจอร์ดังกล่าว ด้วยวิธี เค-เนียร์สเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัด อัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.18 และ 4.19 ตามลำดับ รายละเอียดของพีเจอร์ที่ได้แสดงในภาคผนวก ข

4.5.3 ผลการทดลองของชุดเอกสาร PDDP J-series เมื่อใช้พีเจอร์แบบคำเดียว

สร้างพีเจอร์แบบคำเดียว จากทุกคำในทุกเอกสารในชุดเรียนรู้ แล้วทำการเลือกพีเจอร์ที่มีค่า คะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 2 เอกสาร และทำการ ทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ ดัง แสดงในตารางที่ 4.20 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วย พีเจอร์ดังกล่าว ด้วยวิธี เค-เนียร์สเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัด อัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.21 และ 4.22 ตามลำดับ

ตารางที่ 4.16 แสดงจำนวนเอกสาร J-series ที่ใช้ในการเรียนรู้และทดสอบในแต่ละชุดการทดลอง

ชื่อประเภท (ตัวย่อ)	จำนวนเอกสารในชุดการทดลอง (fold) ที่							
	1		2		3		4	
	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ
A	15	5	15	5	15	5	15	5
B	14	5	14	5	14	5	15	4
C	14	5	14	5	15	4	14	5
E	14	5	15	4	14	5	14	5
I	15	4	14	5	14	5	14	5
L	12	4	12	4	12	4	12	4
M	13	4	13	4	13	4	12	5
P	14	5	14	5	14	5	15	4
S	14	4	14	4	13	5	13	5
Z	14	5	14	5	14	5	15	4
รวม	139	46	139	46	138	47	139	46
	185		185		185		185	

ตารางที่ 4.17 แสดงจำนวนฟีเจอร์จากต้นไม้ลิ ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง
ของชุดเอกสาร PDDP J-series

ชุดที่	จำนวนฟีเจอร์จากต้นไม้ลิ	
	ทั้งหมด	ฟีเจอร์ที่ถูกเลือก
1	70,248	2,747
2	63,872	2,543
3	66,730	2,552
4	66,086	2,455
ค่าเฉลี่ย	66,734	2,574

ตารางที่ 4.18 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์ ค่า $k = 1$
ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์จากต้นไม้ลิ ของชุดเอกสาร PDDP J-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	69.57	73.91	71.74	71.74	71.74	76.09	73.91	76.09	69.57	76.09
2	78.26	73.91	84.78	86.96	84.78	89.13	86.96	82.61	80.44	89.13
3	76.60	80.85	78.72	85.11	87.23	89.36	91.49	91.49	91.49	91.49
4	73.91	71.74	78.26	78.26	80.44	80.44	84.78	82.61	86.96	86.96
ค่าเฉลี่ย	74.58	75.10	78.38	80.52	81.05	83.75	84.29	83.20	82.11	85.92

ตารางที่ 4.19 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์
ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์จากต้นไม้ลิ ของชุดเอกสาร J-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.7015	0.7268	0.7005	0.7117	0.7167	0.7629	0.7358	0.7690	0.6742	0.7690
2	0.7416	0.7223	0.8409	0.8712	0.8460	0.8902	0.8685	0.8204	0.7932	0.8902
3	0.7575	0.7981	0.7864	0.8437	0.8676	0.8827	0.9134	0.9134	0.9134	0.9134
4	0.7342	0.6973	0.7856	0.7822	0.8001	0.7982	0.8457	0.8230	0.8698	0.8698
ค่าเฉลี่ย	0.7337	0.7361	0.7784	0.8022	0.8076	0.8335	0.8409	0.8315	0.8127	0.8606

ตารางที่ 4.20 แสดงจำนวนฟีเจอร์แบบคำเดี่ยวที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง
ของชุดเอกสาร PDDP J-series

ชุดที่	จำนวนฟีเจอร์แบบคำเดี่ยว	
	ทั้งหมด	ฟีเจอร์ที่ถูกเลือก
1	12,568	3,835
2	11,799	3,588
3	11,629	3,532
4	11,164	3,476
ค่าเฉลี่ย	11,790	3,608

ตารางที่ 4.21 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า $k = 1$
ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์แบบคำเดี่ยว ของชุดเอกสาร PDDP J-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	65.22	73.91	71.74	73.91	73.91	69.57	71.74	69.57	71.74	73.91
2	80.44	73.91	82.61	78.26	80.44	82.61	80.44	76.09	76.09	82.61
3	74.47	82.98	82.98	85.11	82.98	87.23	85.11	85.11	85.11	87.23
4	73.91	76.09	84.78	71.74	76.09	78.26	80.44	78.26	80.44	84.78
ค่าเฉลี่ย	73.51	76.72	80.53	77.26	78.35	79.42	79.43	77.26	78.34	82.14

ตารางที่ 4.22 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์
ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์แบบคำเดี่ยว ของชุดเอกสาร J-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.6506	0.7462	0.7069	0.7301	0.7295	0.6941	0.7071	0.6991	0.7138	0.7462
2	0.7888	0.7216	0.8251	0.7811	0.8023	0.8241	0.7958	0.7443	0.7343	0.8251
3	0.7403	0.8241	0.8192	0.8199	0.8002	0.8717	0.8476	0.8476	0.8439	0.8717
4	0.7178	0.7430	0.8492	0.6802	0.7451	0.7670	0.7872	0.7458	0.7872	0.8492
ค่าเฉลี่ย	0.7244	0.7587	0.8001	0.7528	0.7693	0.7892	0.7844	0.7592	0.7698	0.8230

4.5.4 ผลการทดลองของชุดเอกสาร PDDP J-series เมื่อใช้พีเจอร์แบบคำเดี่ยว และไป-แกรม

สร้างพีเจอร์ที่เป็นคำเดี่ยว และไป-แกรม โดยสร้างไป-แกรมจากวลีที่สกัดได้ แล้วทำการเลือกพีเจอร์ที่มีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 2 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.23 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอร์ดังกล่าว ด้วยวิธี เค-เนียร์สเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.24 และ 4.25 ตามลำดับ

ตารางที่ 4.23 แสดงจำนวนคำเดี่ยวและไป-แกรมที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลองของชุดเอกสาร PDDP J-series

ชุดที่	จำนวนพีเจอร์คำเดี่ยวและไป-แกรม	
	ทั้งหมด	พีเจอร์ที่ถูกเลือก
1	49,872	5,632
2	45,610	4,957
3	47,220	4,897
4	46,766	4,782
ค่าเฉลี่ย	47,367	5,067

ตารางที่ 4.24 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า k = 1 ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์แบบคำเดี่ยวและไป-แกรม ของชุดเอกสาร PDDP J-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	65.22	78.26	67.39	69.57	76.09	67.39	71.74	69.57	69.57	78.26
2	80.44	76.09	80.44	82.61	82.61	80.44	80.44	78.26	78.26	82.61
3	74.47	82.98	82.98	80.85	80.85	87.23	87.23	89.36	89.36	89.36
4	76.09	76.09	82.61	71.74	73.91	78.26	78.26	76.09	78.26	82.61
ค่าเฉลี่ย	74.05	78.35	78.35	76.19	78.37	78.33	79.42	78.32	78.86	83.21

ตารางที่ 4.25 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เน็ยรสนเนเบอร์
ค่า $k = 1$ ถึง 9 ในแต่ละชุดการเมื่อใช้พีเจอรแบบคำเดี่ยวและไบ-แกรม ของชุดเอกสาร
J-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.6250	0.7788	0.6555	0.6953	0.7527	0.6656	0.7191	0.6977	0.7011	0.7788
2	0.7888	0.7444	0.7826	0.8078	0.8202	0.7861	0.7950	0.7663	0.7592	0.8202
3	0.7455	0.8241	0.8192	0.7807	0.7755	0.8717	0.8747	0.8960	0.8914	0.8960
4	0.7362	0.7430	0.8268	0.6935	0.7097	0.7516	0.7458	0.7333	0.7629	0.8268
ค่าเฉลี่ย	0.7239	0.7726	0.7710	0.7443	0.7645	0.7688	0.7837	0.7733	0.7787	0.8305

4.5.5 ผลการทดลองของชุดเอกสาร PDDP J-series เมื่อใช้พีเจอรคำเดี่ยวและวลี

สร้างพีเจอรจากวลีที่สกัดได้และคำเดี่ยว แล้วทำการเลือกพีเจอรที่มีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 2 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอรจากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.26 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอรดังกล่าว ด้วยวิธีเค-เน็ยรสนเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.27 และ 4.28 ตามลำดับ

4.5.6 เปรียบเทียบผลการทดลองเมื่อใช้พีเจอรแบบต่าง ของชุดเอกสาร J-series

รูปที่ 4.16 ถึง 4.41 แสดงการเปรียบเทียบผลการทดลองเมื่อใช้พีเจอรแต่ละแบบ โดยเปรียบเทียบทั้งในด้านจำนวนพีเจอร อัตราความถูกต้อง และค่า F-measure

ตารางที่ 4.26 แสดงจำนวนพีเจอรคำเดี่ยวและวลี ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลองของชุดเอกสาร PDDP J-series

ชุดที่	จำนวนพีเจอรคำเดี่ยวและวลี	
	ทั้งหมด	พีเจอรที่ถูกเลือก
1	44,694	5,082
2	40,581	4,504
3	42,333	4,473
4	41,654	4,341
ค่าเฉลี่ย	42,316	4,600

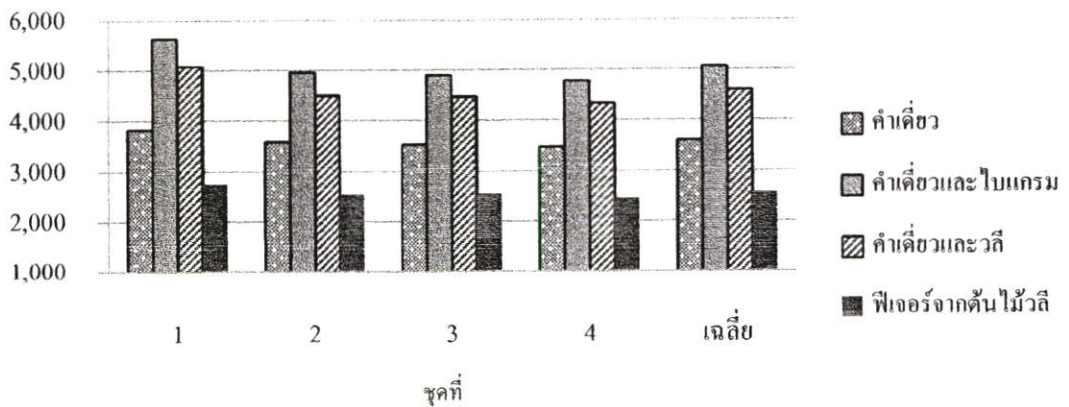
ตารางที่ 4.27 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์ค่าเดียวและวลี ของชุดเอกสาร PDDP J-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	63.04	78.26	67.39	69.57	76.09	67.39	71.74	69.57	73.91	78.26
2	80.44	76.09	84.78	78.26	84.78	82.61	82.61	78.26	78.26	84.78
3	76.60	82.98	82.98	82.98	80.85	87.23	89.36	89.36	89.36	89.36
4	73.91	76.09	84.78	71.74	78.26	76.09	78.26	78.26	80.44	84.78
ค่าเฉลี่ย	73.50	78.35	79.98	75.64	80.00	78.33	80.49	78.86	80.49	84.30

ตารางที่ 4.28 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์ค่าเดียวและวลี ของชุดเอกสาร J-series

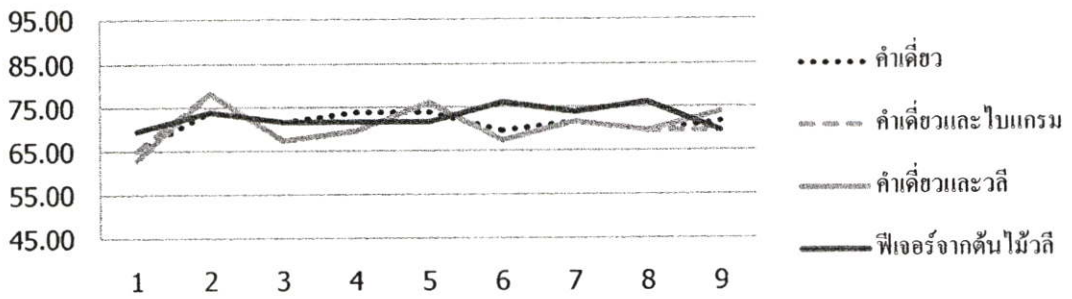
ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.6670	0.7807	0.7802	0.7312	0.7859	0.7943	0.7599	0.7784	0.7623	0.7943
2	0.7930	0.8870	0.8102	0.8669	0.8752	0.8221	0.7663	0.7663	0.7861	0.8870
3	0.7507	0.6915	0.7269	0.8261	0.7766	0.8062	0.7839	0.8927	0.8927	0.8927
4	0.7459	0.7976	0.7886	0.7969	0.7702	0.7702	0.7519	0.7068	0.7654	0.7976
ค่าเฉลี่ย	0.7391	0.7892	0.7765	0.8053	0.8020	0.7982	0.7655	0.7860	0.8016	0.8429

จำนวนฟีเจอร์ที่ใช้ของชุดเอกสาร J-series



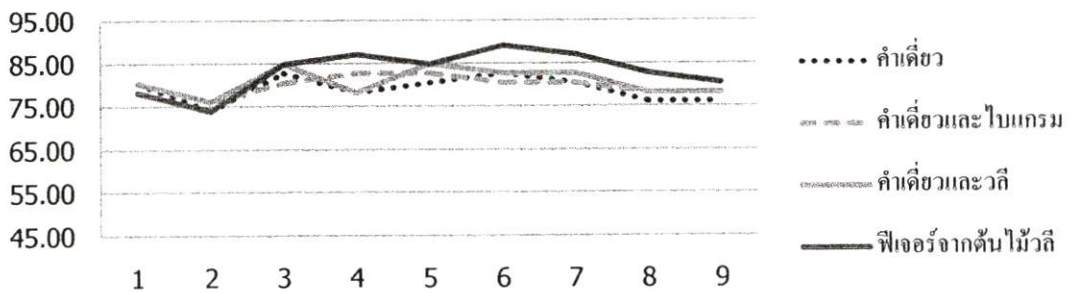
รูปที่ 4.16 แผนภูมิแท่งเปรียบเทียบจำนวนฟีเจอร์ที่ใช้ในแต่ละแบบของชุดเอกสาร J-series

อัตราการถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร J-series ชุดที่ 1



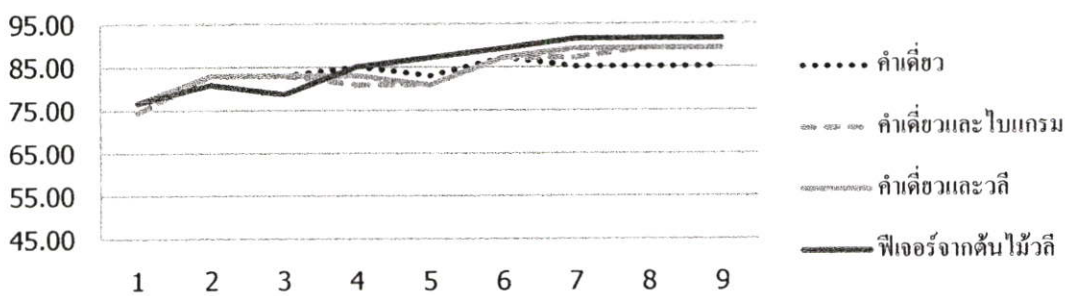
รูปที่ 4.17 กราฟเปรียบเทียบอัตราการถูกต้องในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 1

อัตราการถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร J-series ชุดที่ 2



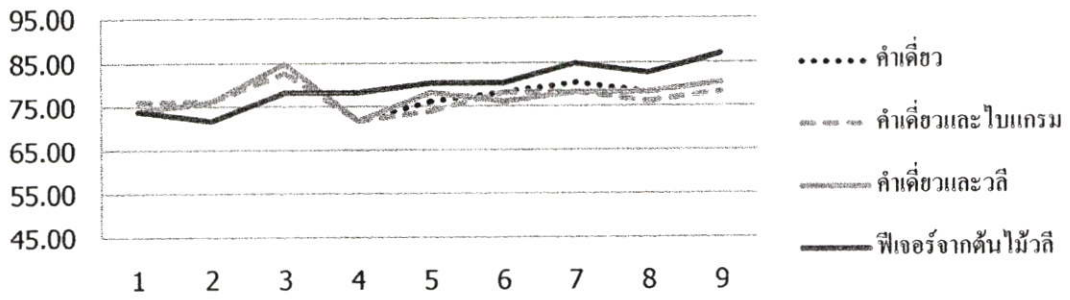
รูปที่ 4.18 กราฟเปรียบเทียบอัตราการถูกต้องในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 2

อัตราการถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร J-series ชุดที่ 3



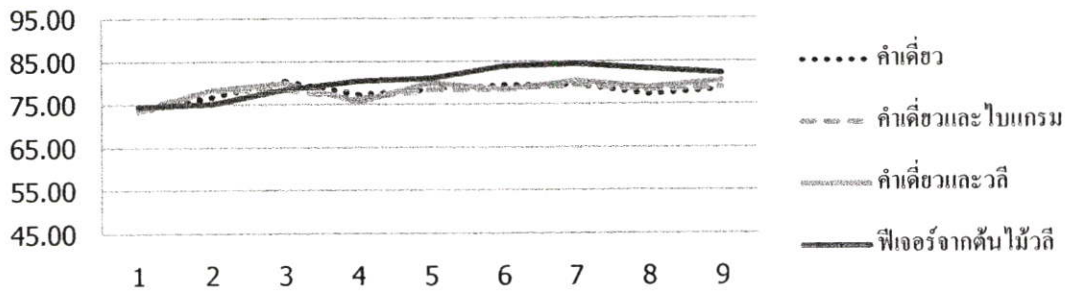
รูปที่ 4.19 กราฟเปรียบเทียบอัตราการถูกต้องในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 3

อัตราการถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร J-series ชุดที่ 4



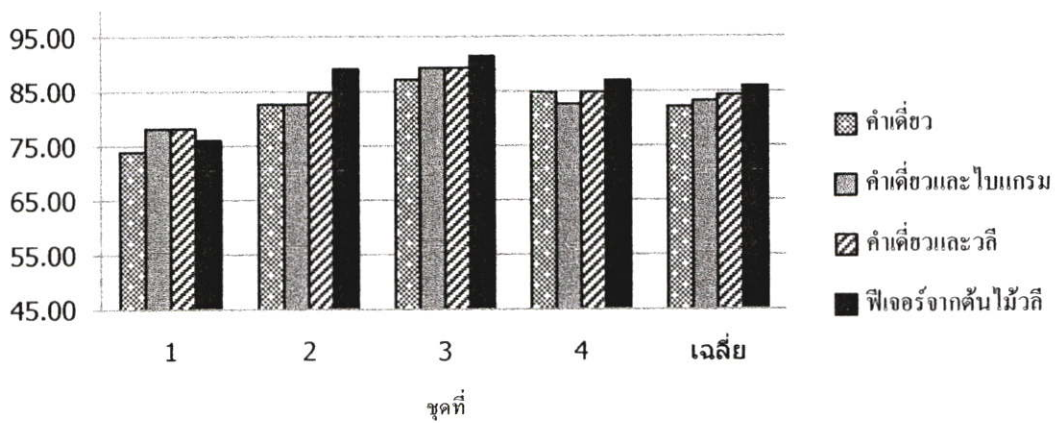
รูปที่ 4.20 กราฟเปรียบเทียบอัตราการถูกต้องในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 4

อัตราการถูกต้องเฉลี่ยเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร J-series



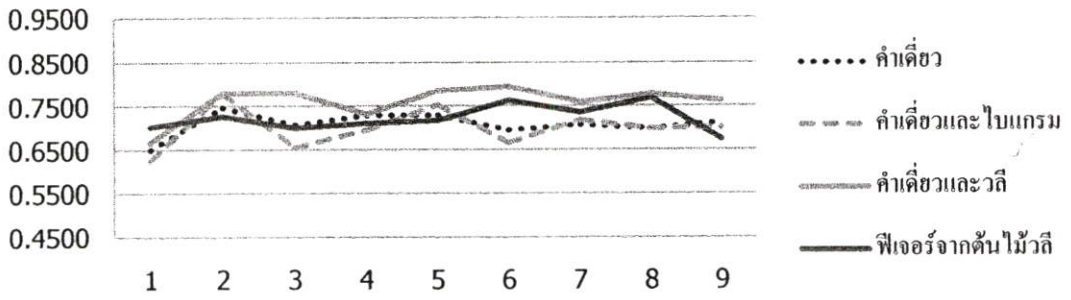
รูปที่ 4.21 กราฟเปรียบเทียบอัตราการถูกต้องเฉลี่ยในแต่ละแบบของชุดเอกสาร J-series

อัตราการถูกต้องสูงสุดของชุดเอกสาร J-series



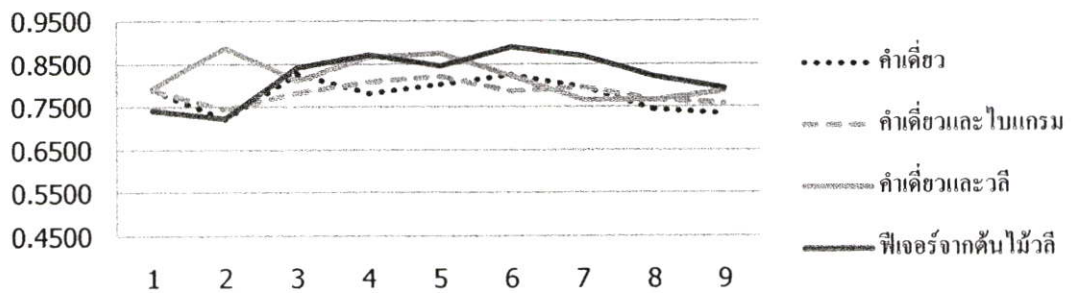
รูปที่ 4.22 แผนภูมิแท่งเปรียบเทียบอัตราการถูกต้องสูงสุดในแต่ละแบบของชุดเอกสาร J-series

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร J-series ชุดที่ 1



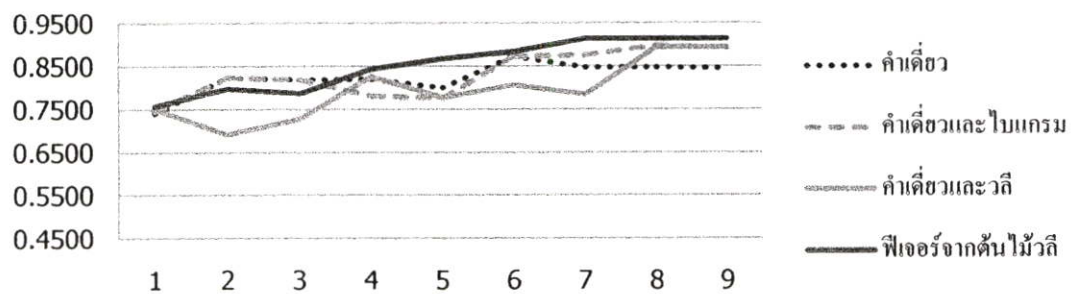
รูปที่ 4.23 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 1

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร J-series ชุดที่ 2



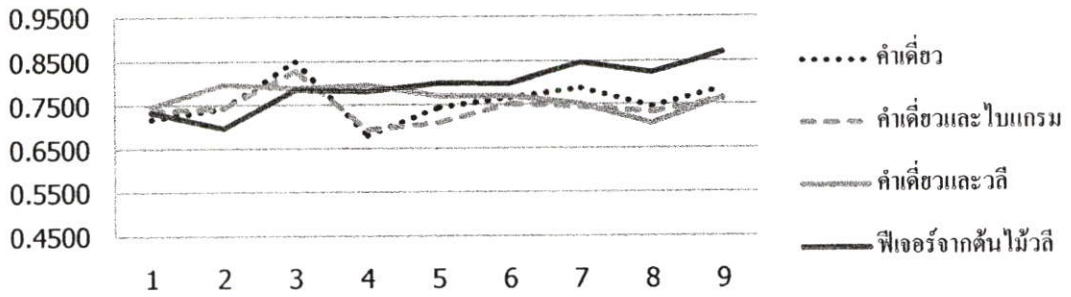
รูปที่ 4.24 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 2

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร J-series ชุดที่ 3



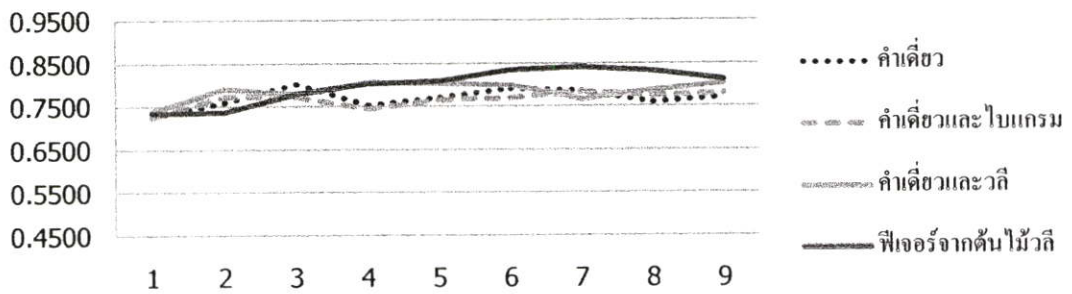
รูปที่ 4.25 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 3

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร J-series ชุดที่ 4



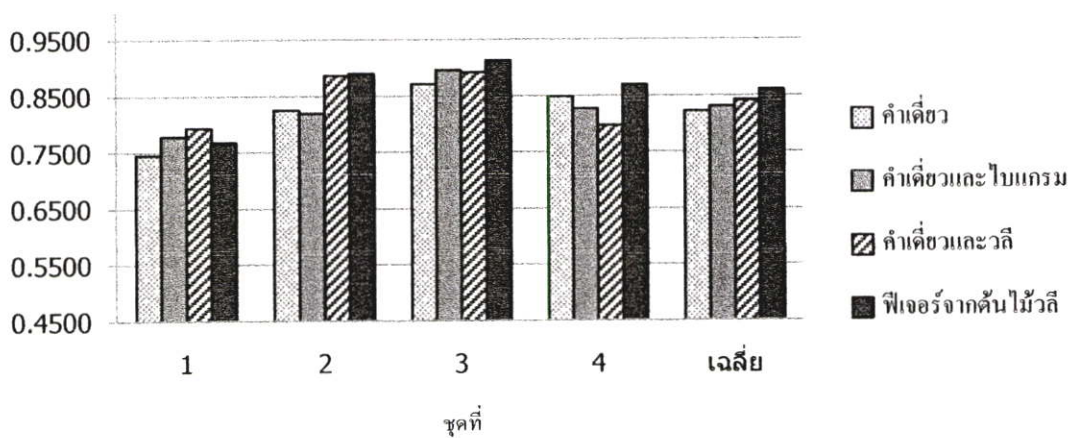
รูปที่ 4.26 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร J-series ชุดที่ 4

ค่า F-measure เฉลี่ยเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร J-series



รูปที่ 4.27 กราฟเปรียบเทียบ ค่า F-measure เฉลี่ยในแต่ละแบบของชุดเอกสาร J-series

ค่า F-measure สูงสุดของชุดเอกสาร J-series



รูปที่ 4.28 แผนภูมิแท่งเปรียบเทียบ ค่า F-measure สูงสุดในแต่ละแบบของชุดเอกสาร J-series

จากรูปที่ 4.16 จะเห็นได้ว่าพีเจอร์ที่ได้จากต้นไม้วลีมีจำนวนน้อยกว่าพีเจอร์แบบอื่นๆ และเมื่อพิจารณาถึงความถูกต้องในการจำแนกประเภทเอกสาร จากรูปที่ 4.22 และ 4.28 พีเจอร์จากต้นไม้วลีให้อัตราความถูกต้องและค่า F-measure เฉลี่ยจากทุกชุดการทดลองสูงกว่าพีเจอร์แบบอื่นๆ ซึ่งเป็นไปตามเป้าหมายของเราที่ต้องการได้พีเจอร์ที่มีประสิทธิภาพ โดยใช้วลีร่วมกับคำเดี่ยว ที่ให้ความถูกต้องในการจำแนกประเภทเอกสารสูงแต่มีจำนวนน้อย

4.6 ชุดเอกสาร PDDP K-series ที่ใช้ในการทดลองและผลการทดลอง

4.6.1 ชุดเอกสาร PDDP K-series ที่ใช้ในการทดลอง

ชุดเอกสาร PDDP K-series มีเอกสารทั้งหมด 2,340 แบ่งเป็น 20 ประเภท ทำการสุ่มตัวอย่างแบบชั้นภูมิ โดยใช้อัตราส่วนในการเลือกเอกสารจากแต่ละชั้นภูมิเท่ากับ 0.1 โดยรายละเอียดในการสุ่มตัวอย่างเอกสารแสดงในภาคผนวก ก โดยในขั้นแรกสุ่มตัวอย่างเอกสารได้เอกสารทั้งหมด 796 เอกสาร และจากนั้นเลือกประเภทเอกสารที่มีจำนวนเอกสารน้อยออก คือ ประเภทเอกสาร Entertainment, Art, Cable, Media, Multimedia, Stage และ Technology ทำให้เหลือประเภทเอกสารทั้งหมด 13 ประเภท และเอกสารตัวอย่าง 387 เอกสาร

ในการทดลองจำแนกประเภทเอกสารได้แบ่งชุดเอกสาร PDDP K-series เพื่อใช้ในการเรียนรู้ และใช้ในการทดสอบ สำหรับทำ 4-fold cross validation โดยการสุ่มเลือกเอกสารได้จำนวนเอกสารสำหรับเรียนรู้และทดสอบในแต่ละชุดการทดลองดังแสดงในตาราง 4.29 และทดลองเพื่อวัดประสิทธิภาพของพีเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลีในการจำแนกประเภทเอกสารเปรียบเทียบกับพีเจอร์แบบอื่นๆ โดยทำการทดลองสร้างพีเจอร์ 4 แบบ สำหรับการจำแนกประเภทคือ 1.พีเจอร์จากต้นไม้วลี 2.คำเดี่ยว 3.คำเดี่ยวและไบนารี-แกรม และ 4.คำเดี่ยวและวลี

4.6.2 ผลการทดลองของชุดเอกสาร PDDP K-series เมื่อใช้พีเจอร์ที่ได้จากการเลือกพีเจอร์โดยใช้ต้นไม้วลี

สร้างต้นไม้วลีและทำการเลือกพีเจอร์โดยใช้ต้นไม้วลีโดยกำหนดให้แต่ละเทอมมีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 2 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ดังแสดงในตารางที่ 4.30 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอร์ดังกล่าว ด้วยวิธี เค-เน็ยเรสเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.31 และ 4.32 ตามลำดับรายละเอียดของพีเจอร์ที่ได้แสดงในภาคผนวก ข

4.6.3 ผลการทดลองของชุดเอกสาร PDDP K-series เมื่อใช้พีเจอร์แบบคำเดียว

สร้างพีเจอร์แบบคำเดียว จากทุกคำในทุกเอกสารในชุดเรียนรู้ แล้วทำการเลือกพีเจอร์ที่มีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 2 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.33 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอร์ดังกล่าว ด้วยวิธี เค-เนียร์สเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.34 และ 4.35 ตามลำดับ

ตารางที่ 4.29 แสดงจำนวนเอกสาร K-series ที่ใช้ในการเรียนรู้และทดสอบในแต่ละชุดการทดลอง

ชื่อประเภท	จำนวนเอกสารในชุดการทดลอง (fold) ที่							
	1		2		3		4	
	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ
Business	13	5	13	5	14	4	14	4
Culture	9	3	9	3	9	3	9	3
Film	39	13	39	13	39	13	39	13
Industry	7	2	7	2	6	3	7	2
Music	16	5	16	5	15	6	16	5
Online	6	2	6	2	6	2	6	2
People	36	12	36	12	36	12	36	12
Review	22	8	22	8	23	7	23	7
Television	25	8	25	8	24	9	25	8
Variety	7	3	8	2	7	3	8	2
Health	85	28	84	29	85	28	85	28
Politics	11	3	11	3	10	4	10	4
Sports	15	4	14	5	14	5	14	5
รวม	291	96	290	97	288	99	292	95
	387		387		387		387	

ตารางที่ 4.30 แสดงจำนวนฟีเจอร์จากต้นไม้ที่ ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง
ของชุดเอกสาร PDDP K-series

ชุดที่	จำนวนฟีเจอร์จากต้นไม้	
	ทั้งหมด	ฟีเจอร์ที่ถูกเลือก
1	40,802	2,685
2	41,616	2,745
3	41,945	2,721
4	42,597	2,900
ค่าเฉลี่ย	41,740	2,763

ตารางที่ 4.31 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์ ค่า $k = 1$
ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์จากต้นไม้ ของชุดเอกสาร PDDP K-
series

ชุด ที่	k									ค่า สูงสุด
	1	2	3	4	5	6	7	8	9	
1	79.17	77.08	76.04	77.08	78.13	83.33	81.25	83.33	86.46	86.46
2	72.17	73.20	75.26	75.26	76.29	76.29	75.26	79.38	81.44	81.44
3	77.32	73.20	74.23	74.23	72.17	74.23	77.32	74.23	74.23	77.32
4	76.29	81.44	82.47	84.54	86.60	85.57	85.57	84.54	84.54	86.60
ค่าเฉลี่ย	76.24	76.23	77.00	77.78	78.29	79.85	79.85	80.37	81.67	82.96

ตารางที่ 4.32 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์
ค่า $k = 1$ ถึง 9 เมื่อใช้ฟีเจอร์จากต้นไม้ ของชุดเอกสาร K-series

ชุด ที่	k									ค่า สูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.7938	0.7694	0.7604	0.7723	0.7749	0.8168	0.7959	0.8164	0.8499	0.8499
2	0.7160	0.7180	0.7564	0.7515	0.7575	0.7545	0.7464	0.7908	0.8132	0.8132
3	0.7715	0.7339	0.7467	0.7424	0.7273	0.7489	0.7734	0.7374	0.7422	0.7734
4	0.7519	0.8162	0.8267	0.8480	0.8692	0.8611	0.8611	0.8396	0.8394	0.8692
ค่าเฉลี่ย	0.7583	0.7594	0.7725	0.7786	0.7822	0.7953	0.7942	0.7961	0.8112	0.8264

ตารางที่ 4.33 แสดงจำนวนฟีเจอร์แบบคำเดี่ยวที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง
ของชุดเอกสาร PDDP K-series

ชุดที่	จำนวนฟีเจอร์แบบคำเดี่ยว	
	ทั้งหมด	ฟีเจอร์ที่ถูกเลือก
1	10,520	3,672
2	10,650	3,618
3	10,867	3,777
4	10,872	3,803
ค่าเฉลี่ย	10,727	3,718

ตารางที่ 4.34 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์ ค่า $k = 1$
ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์แบบคำเดี่ยว ของชุดเอกสาร PDDP K-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	76.04	77.08	73.96	79.17	83.33	84.38	83.33	82.29	83.33	84.38
2	64.95	70.10	72.17	70.10	71.13	74.23	72.17	73.20	75.26	75.26
3	71.13	76.29	79.38	74.23	76.29	80.41	78.35	76.29	75.26	80.41
4	67.01	71.13	76.29	80.41	79.38	81.44	80.41	83.51	82.47	83.51
ค่าเฉลี่ย	69.78	73.65	75.45	75.98	77.53	80.11	78.57	78.82	79.08	80.89

ตารางที่ 4.35 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์
ค่า $k = 1$ ถึง 9 เมื่อใช้ฟีเจอร์แบบคำเดี่ยว ของชุดเอกสาร K-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.7655	0.7595	0.7352	0.7762	0.8156	0.8237	0.8114	0.8018	0.8192	0.8237
2	0.6513	0.6885	0.7145	0.7036	0.7159	0.7415	0.7123	0.7245	0.7372	0.7415
3	0.7032	0.7611	0.7913	0.7418	0.7679	0.8005	0.7833	0.7641	0.7510	0.8005
4	0.6528	0.7114	0.7709	0.8112	0.8067	0.8174	0.8176	0.8404	0.8229	0.8404
ค่าเฉลี่ย	0.6932	0.7301	0.7530	0.7582	0.7765	0.7958	0.7812	0.7827	0.7825	0.8015

4.6.4 ผลการทดลองของชุดเอกสาร PDDP K-series เมื่อใช้พีเจอร์แบบคำเดี่ยว และไบ-แกรม

สร้างพีเจอร์ที่เป็นคำเดี่ยว และไบ-แกรม โดยสร้างไบ-แกรมจากวลีที่สกัดได้ แล้วทำการเลือกพีเจอร์ที่มีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 2 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.36 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอร์ดังกล่าว ด้วยวิธี เค-เน็ยเรสเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.37 และ 4.38 ตามลำดับ

ตารางที่ 4.36 แสดงจำนวนคำเดี่ยวและไบ-แกรมที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลองของชุดเอกสาร PDDP K-series

ชุดที่	จำนวนพีเจอร์คำเดี่ยวและไบ-แกรม	
	ทั้งหมด	พีเจอร์ที่ถูกเลือก
1	29,870	4,497
2	30,520	4,430
3	30,902	4,575
4	31,281	4,661
ค่าเฉลี่ย	30,643	4,541

ตารางที่ 4.37 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เน็ยเรสเนเบอร์ ค่า k = 1 ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์แบบคำเดี่ยวและไบ-แกรม ของชุดเอกสาร PDDP K-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	76.04	77.08	75.00	80.21	82.29	83.33	84.38	83.33	83.33	84.38
2	67.01	73.20	72.17	71.13	74.23	73.20	72.17	75.26	73.20	75.26
3	70.10	78.35	78.35	76.29	76.29	80.41	80.41	78.35	77.32	80.41
4	64.95	73.20	80.41	83.51	81.44	82.47	78.35	83.51	83.51	83.51
ค่าเฉลี่ย	69.53	75.46	76.48	77.78	78.56	79.85	78.83	80.11	79.34	80.89

ตารางที่ 4.38 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เน็ยเรสเนเบอร์
ค่า $k = 1$ ถึง 9 ในแต่ละชุดการเมื่อใช้พีเจอร์แบบคำเดี่ยวและไบ-แกรม ของชุดเอกสาร
K-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.7633	0.7685	0.7402	0.7857	0.8062	0.8135	0.8226	0.8117	0.8193	0.8226
2	0.6739	0.7155	0.7208	0.7135	0.7440	0.7203	0.7114	0.7428	0.7200	0.7440
3	0.6978	0.7844	0.7810	0.7644	0.7661	0.8065	0.8053	0.7852	0.7730	0.8065
4	0.6325	0.7313	0.8053	0.8391	0.8243	0.8291	0.7970	0.8468	0.8452	0.8468
ค่าเฉลี่ย	0.6919	0.7499	0.7619	0.7757	0.7851	0.7923	0.7841	0.7966	0.7893	0.8050

4.6.5 ผลการทดลองของชุดเอกสาร PDDP K-series เมื่อใช้พีเจอร์คำเดี่ยวและวลี

สร้างพีเจอร์จากวลีที่สกัดได้และคำเดี่ยว แล้วทำการเลือกพีเจอร์ที่มีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 2 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.39 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอร์ดังกล่าว ด้วยวิธีเค-เน็ยเรสเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.40 และ 4.41 ตามลำดับ

ตารางที่ 4.39 แสดงจำนวนพีเจอร์คำเดี่ยวและวลี ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง
ของชุดเอกสาร PDDP K-serie

ชุดที่	จำนวนพีเจอร์คำเดี่ยวและวลี	
	ทั้งหมด	พีเจอร์ที่ถูกเลือก
1	25,063	4,201
2	25,607	4,177
3	25,894	4,280
4	26,213	4,390
ค่าเฉลี่ย	25,694	4,262

ตารางที่ 4.40 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์ค่าเดียวและวลี ของชุดเอกสาร PDDP K-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	77.08	73.96	77.08	81.25	82.29	80.21	80.21	81.25	79.17	82.29
2	68.04	68.04	72.17	77.32	76.29	79.38	79.38	80.41	80.41	80.41
3	63.92	73.20	72.17	72.17	76.29	73.20	76.29	74.23	74.23	76.29
4	64.95	73.20	79.38	80.41	81.44	77.32	81.44	81.44	82.47	82.47
ค่าเฉลี่ย	68.50	72.10	75.20	77.79	79.08	77.53	79.33	79.33	79.07	80.37

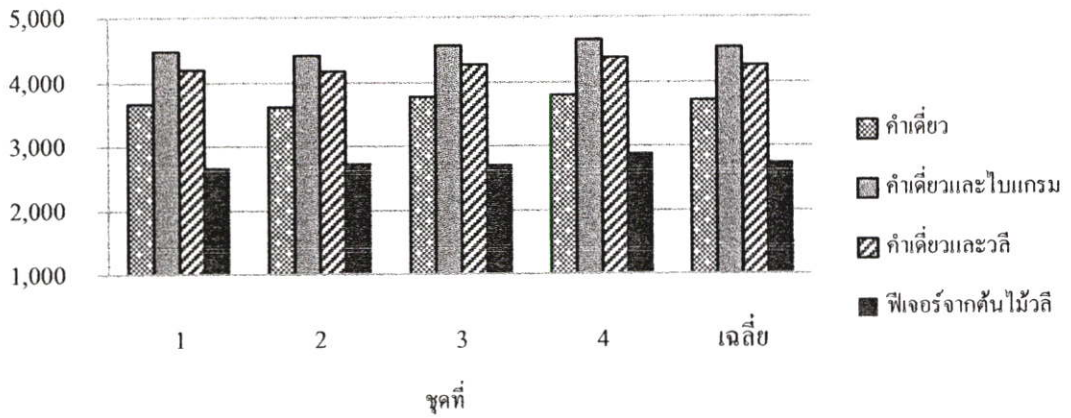
ตารางที่ 4.41 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์ค่าเดียวและวลี ของชุดเอกสาร K-series

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.7633	0.7300	0.7609	0.7973	0.8115	0.7840	0.7850	0.7970	0.7794	0.8115
2	0.6892	0.6718	0.7270	0.7694	0.7680	0.7921	0.7891	0.7974	0.7950	0.7974
3	0.6362	0.7228	0.7235	0.7108	0.7593	0.7303	0.7600	0.7368	0.7441	0.7600
4	0.6364	0.7138	0.7932	0.8028	0.8137	0.7707	0.8100	0.8139	0.8235	0.8235
ค่าเฉลี่ย	0.6813	0.7096	0.7512	0.7701	0.7881	0.7693	0.7860	0.7863	0.7855	0.7981

4.6.6 เปรียบเทียบผลการทดลองเมื่อใช้พีเจอร์แบบต่าง ของชุดเอกสาร K-series

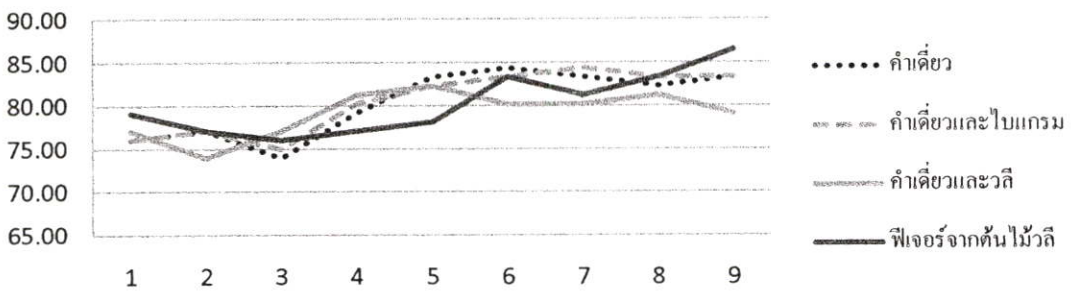
รูปที่ 4.29 ถึง 4.41 แสดงการเปรียบเทียบผลการทดลองเมื่อใช้พีเจอร์แต่ละแบบ โดยเปรียบเทียบทั้งในด้านจำนวนพีเจอร์ อัตราความถูกต้อง และค่า F-measure

จำนวนพีเจอร์ที่ใช้ของชุดเอกสาร K-series



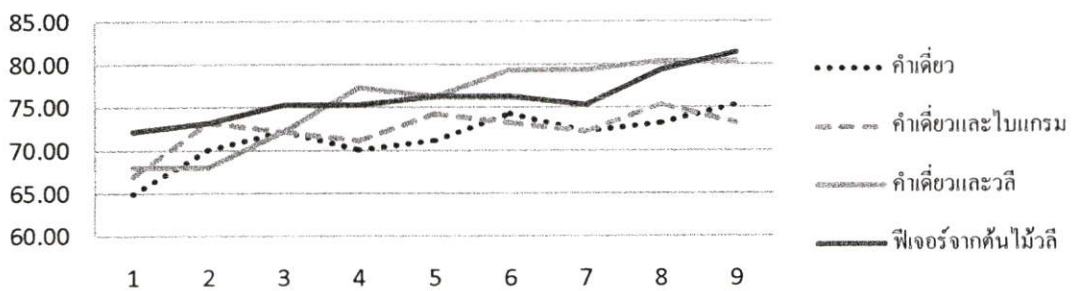
รูปที่ 4.29 แผนภูมิแท่งเปรียบเทียบจำนวนพีเจอร์ที่ใช้ในแต่ละแบบของชุดเอกสาร K-series

อัตราความถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร K-series ชุดที่ 1



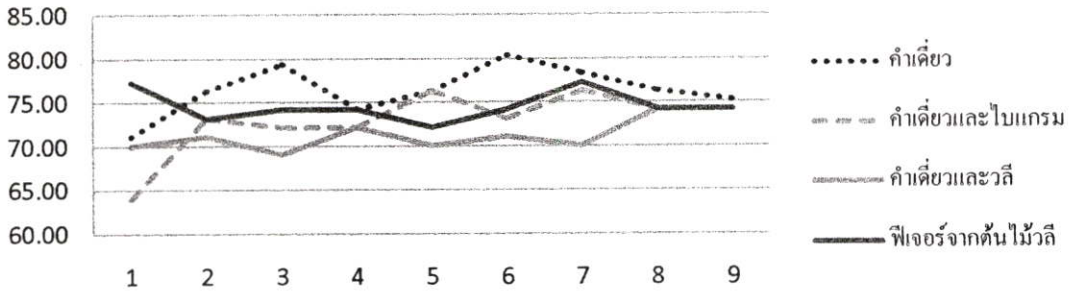
รูปที่ 4.30 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 1

อัตราความถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร K-series ชุดที่ 2



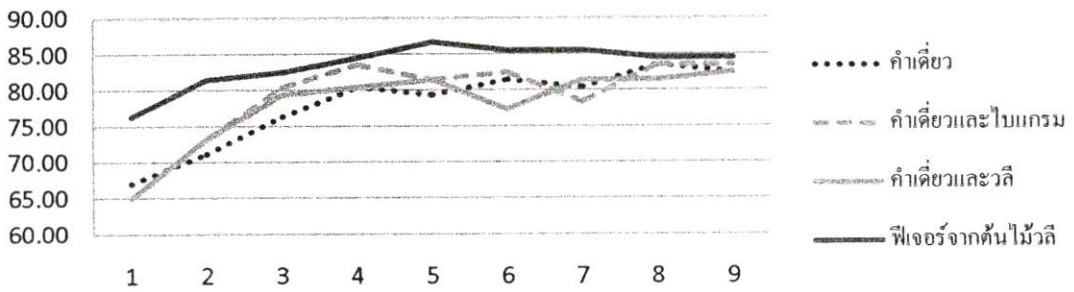
รูปที่ 4.31 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 2

อัตราการถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร K-series ชุดที่ 3



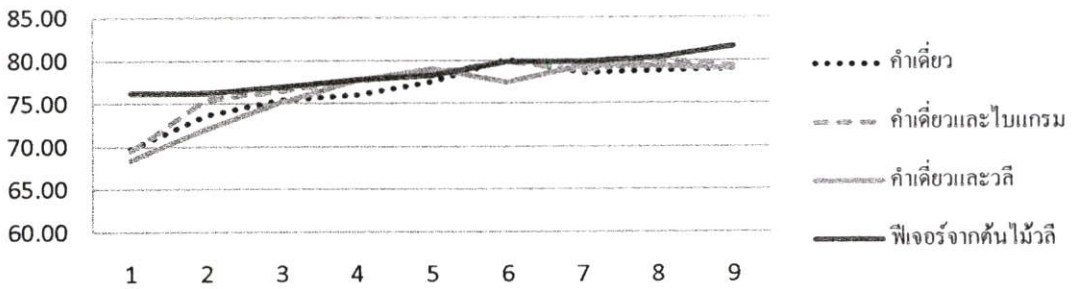
รูปที่ 4.32 กราฟเปรียบเทียบอัตราการถูกต้องในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 3

อัตราการถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร K-series ชุดที่ 4



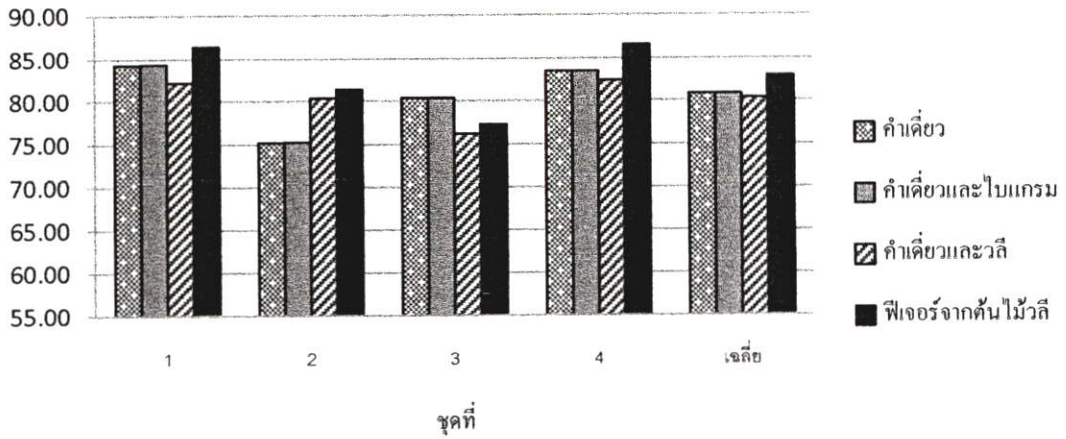
รูปที่ 4.33 กราฟเปรียบเทียบอัตราการถูกต้องในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 4

อัตราการถูกต้องเฉลี่ยเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร K-series



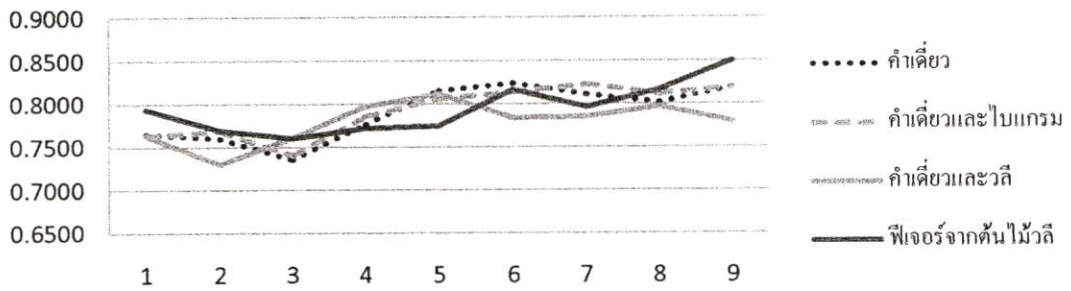
รูปที่ 4.34 กราฟเปรียบเทียบอัตราการถูกต้องเฉลี่ยในแต่ละแบบของชุดเอกสาร K-series

อัตราความถูกต้องสูงสุดของชุดเอกสาร K-series



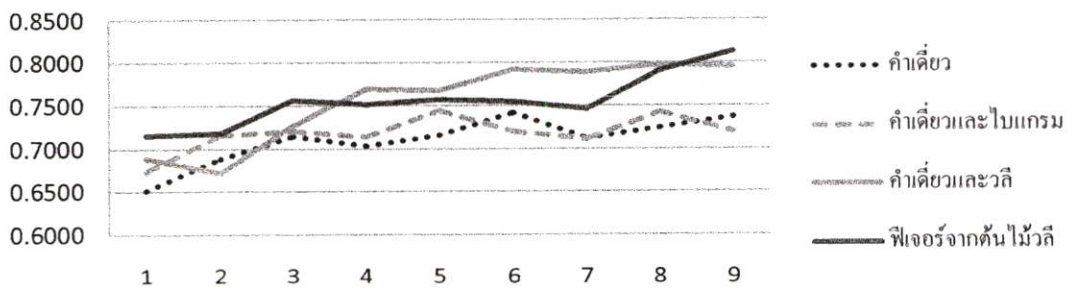
รูปที่ 4.35 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุดในแต่ละแบบของชุดเอกสาร K-series

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร K-series ชุดที่ 1



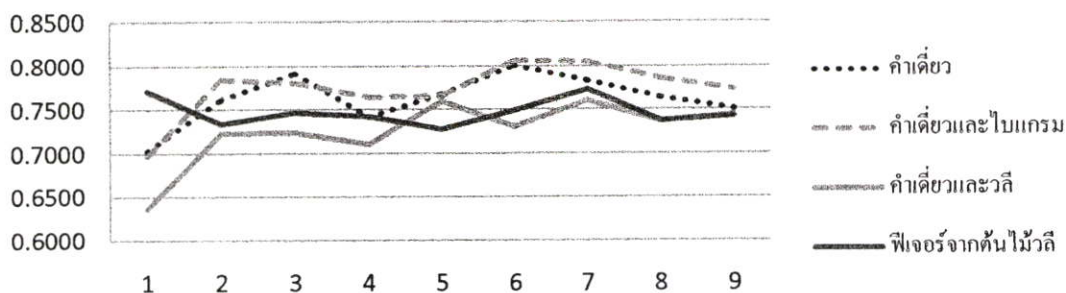
รูปที่ 4.36 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 1

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร K-series ชุดที่ 2



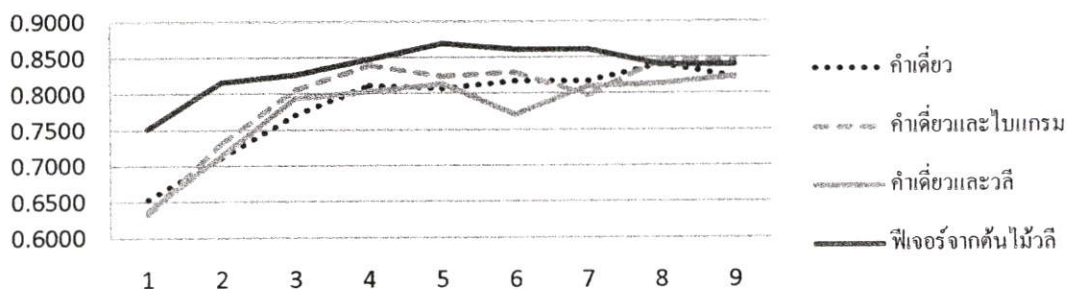
รูปที่ 4.37 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 2

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร K-series ชุดที่ 3



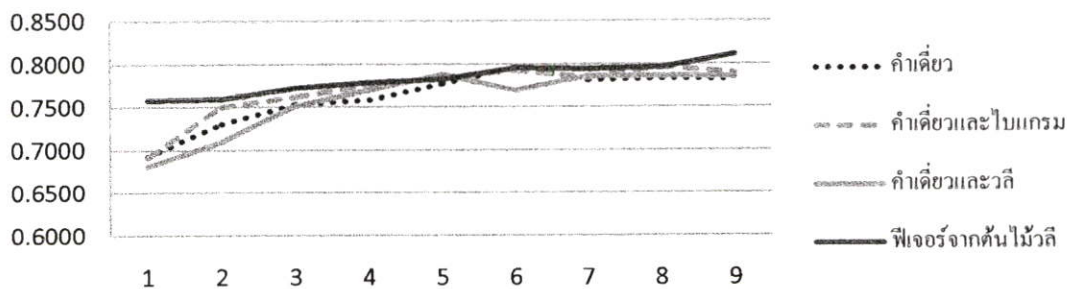
รูปที่ 4.38 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 3

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร K-series ชุดที่ 4



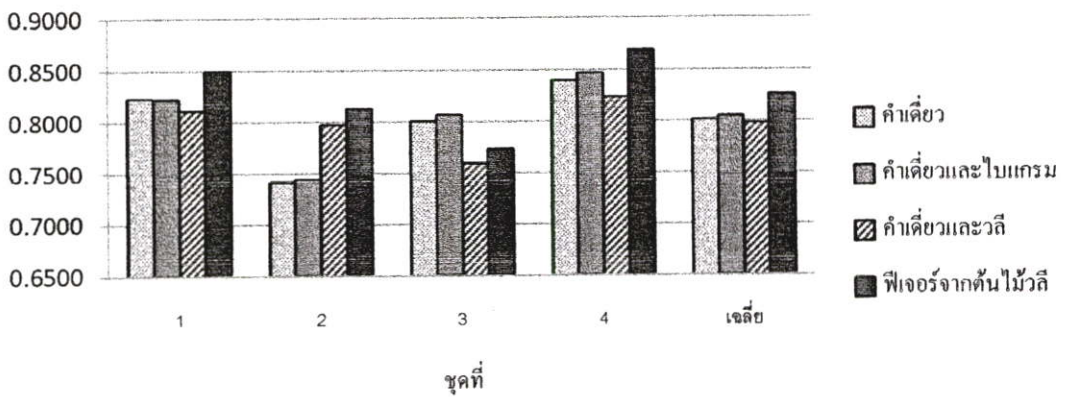
รูปที่ 4.39 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร K-series ชุดที่ 4

ค่า F-measure เฉลี่ยเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร K-series



รูปที่ 4.40 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร K-series

ค่า F-measure สูงสุดของชุดเอกสาร K-series



รูปที่ 4.41 แผนภูมิแท่งเปรียบเทียบ ค่า F-measure สูงสุดในแต่ละแบบของชุดเอกสาร K-series

จากรูปที่ 4.29 จะเห็นได้ว่าฟีเจอร์ที่ได้จากต้นไม้วลีมีจำนวนน้อยกว่าฟีเจอร์แบบอื่นๆ และเมื่อพิจารณาถึงความถูกต้องในการจำแนกประเภทเอกสาร จากรูปที่ 4.35 และ 4.41 ในทุกชุดการทดลอง ฟีเจอร์จากต้นไม้วลีให้อัตราความถูกต้องและค่า F-measure เฉลี่ยสูงกว่าฟีเจอร์แบบอื่นๆ ซึ่งเป็นไปตามเป้าหมายของเราที่ต้องการได้ฟีเจอร์ที่มีประสิทธิภาพ โดยใช้วลีร่วมกับคำเดี่ยว ที่ให้ความถูกต้องในการจำแนกประเภทเอกสารสูงแต่มีจำนวนน้อย

4.7 ชุดเอกสาร Reuters-top10ที่ใช้ในการทดลองและผลการทดลอง

4.7.1 ชุดเอกสาร Reuters-top10ที่ใช้ในการทดลอง

ชุดข้อมูล Reuters-Top10 คัดเลือกมาจาก ชุดข้อมูลข่าว Reuters-21578 [18] ซึ่ง David D. Lewis และคณะได้รวบรวมไว้ ซึ่งมีประเภทเอกสารอยู่ 135 ประเภทเอกสาร จำนวน 21,578 เอกสารข่าว โดย Reuters-Top10 คัดเลือกเอาเฉพาะเอกสารข่าว 10 กลุ่มแรกที่มีจำนวนข่าวมากที่สุดจาก 135 กลุ่มเอกสารข่าว ซึ่งชุดข้อมูล Reuters-Top10 [19] เป็นชุดข้อมูลที่ใช้ในการทดสอบการจำแนกประเภทเอกสารและการจัดกลุ่มเอกสาร ที่ใช้กันอย่างแพร่หลาย ลักษณะการเก็บอยู่ในรูปแบบเพิ่มข้อมูล XML มี 10 กลุ่มเอกสาร จำนวน 2,770 เอกสาร ทำการสุ่มตัวอย่างแบบชั้นภูมิ โดยใช้อัตราส่วนในการเลือกเอกสารจากแต่ละชั้นภูมิเท่ากับ 0.1 โดยรายละเอียดในการสุ่มตัวอย่างเอกสารแสดงในภาคผนวก ก โดยได้เอกสารตัวอย่างทั้งหมด 504

ในการทดลองจำแนกประเภทเอกสารได้แบ่งชุดเอกสาร Reuters-top10 เพื่อใช้ในการเรียนรู้ และใช้ในการทดสอบ สำหรับทำ 4-fold cross validation โดยการสุ่มเลือกเอกสารได้จำนวนเอกสารสำหรับเรียนรู้และทดสอบในแต่ละชุดการทดลองดังแสดงในตาราง 4.42 และทดลองเพื่อวัดประสิทธิภาพของฟีเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลีในการจำแนกประเภทเอกสาร

เปรียบเทียบกับพีเจอร์แบบอื่นๆ โดยทำการทดลองสร้างพีเจอร์ 4 แบบ สำหรับการจำแนกประเภท คือ 1.พีเจอร์จากต้นไม้ตัดสินใจ 2.คำเดียว 3.คำเดียวและไบนารี-แกรม และ 4.คำเดียวและวลี

ตารางที่ 4.42 แสดงจำนวนเอกสาร Reuters-top10ที่ใช้ในการเรียนรู้และทดสอบในแต่ละชุดการทดลอง

ชื่อประเภท	จำนวนเอกสารในชุดการทดลอง (fold) ที่							
	1		2		3		4	
	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ	เรียนรู้	ทดสอบ
acq (a)	102	33	101	34	101	34	101	34
corn (c)	9	3	9	3	9	3	9	3
crude (d)	28	10	28	10	29	9	29	9
earn (e)	127	42	127	42	127	42	126	43
grain (g)	24	7	23	8	23	8	23	8
interest (i)	19	6	19	6	18	7	19	6
money-fx (m)	27	10	28	9	28	9	28	9
ship (s)	14	5	15	4	14	5	14	5
trade (t)	18	6	18	6	18	6	18	6
wheat (w)	11	3	10	4	10	4	11	3
รวม	379	125	378	126	377	127	378	126
	504		504		504		504	

4.7.2 ผลการทดลองของชุดเอกสาร Reuters-top10เมื่อใช้พีเจอร์ที่ได้จากการเลือกพีเจอร์โดยใช้ต้นไม้ตัดสินใจ

สร้างต้นไม้ตัดสินใจและทำการเลือกพีเจอร์โดยใช้ต้นไม้ตัดสินใจโดยกำหนดให้แต่ละเทอมมีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 1 เอกสาร และทำการทดลองจำแนกประเภทเอกสารโดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.43 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอร์ดังกล่าว ด้วยวิธี เค-เนียร์สเนบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.44 และ 4.45 ตามลำดับรายละเอียดของพีเจอร์ที่ได้แสดงในภาคผนวก ข

4.7.3 ผลการทดลองของชุดเอกสาร Reuters-top10 เมื่อใช้พีเจอร์แบบคำเดี่ยว

สร้างพีเจอร์แบบคำเดี่ยว จากทุกคำในทุกเอกสาร ในชุดเรียนรู้ แล้วทำการเลือกพีเจอร์ที่มีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 1 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.46 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอร์ดังกล่าว ด้วยวิธี เค-เนียร์สเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.47 และ 4.48 ตามลำดับ

ตารางที่ 4.43 แสดงจำนวนพีเจอร์จากต้นไม้วลี ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลองของชุดเอกสาร Reuters-top10

ชุดที่	จำนวนพีเจอร์จากต้นไม้วลี	
	ทั้งหมด	พีเจอร์ที่ถูกเลือก
1	15,781	2,095
2	15,788	2,036
3	15,610	2,072
4	15,607	2,195
ค่าเฉลี่ย	15,638	2,100

ตารางที่ 4.44 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียร์สเนเบอร์ ค่า k = 1 ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์จากต้นไม้วลี ของชุดเอกสาร Reuters-top10

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	70.40	76.80	72.80	77.60	80.00	79.20	80.00	76.80	80.00	80.00
2	77.78	82.54	80.16	80.16	81.75	82.54	83.33	82.54	84.92	84.92
3	74.02	76.38	74.80	77.17	79.53	79.53	81.10	79.53	78.74	81.10
4	72.22	71.43	74.60	73.81	73.81	76.19	75.40	76.19	76.19	76.19
ค่าเฉลี่ย	73.60	76.79	75.59	77.18	78.77	79.36	79.96	78.76	79.96	80.55

ตารางที่ 4.45 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิเค-เนียร์เซนเบอร์
ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์จากต้นไม้วลี ของชุดเอกสาร
Reuters-top10

ชุด ที่	k									ค่า สูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.7013	0.7597	0.7226	0.7619	0.7840	0.7826	0.7804	0.7498	0.7877	0.7877
2	0.7764	0.8111	0.7941	0.7984	0.8149	0.8084	0.8226	0.8140	0.8353	0.8353
3	0.7336	0.7555	0.7344	0.7651	0.7887	0.7848	0.8041	0.7886	0.7731	0.8041
4	0.7115	0.7006	0.7408	0.7285	0.7244	0.7556	0.7471	0.7545	0.7520	0.7556
ค่าเฉลี่ย	0.7307	0.7567	0.7480	0.7635	0.7780	0.7828	0.7885	0.7767	0.7870	0.7956

ตารางที่ 4.46 แสดงจำนวนฟีเจอร์แบบคำเดี่ยวที่ได้จากชุดเอกสารเรียนรู้อ ในแต่ละชุดการทดลอง
ของชุดเอกสาร Reuters-top10

ชุดที่	จำนวนฟีเจอร์แบบคำเดี่ยว	
	ทั้งหมด	ฟีเจอร์ที่ถูกเลือก
1	4,103	1,468
2	4,135	1,459
3	4,114	1,438
4	4,084	1,498
ค่าเฉลี่ย	4,109	1,466

ตารางที่ 4.47 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิเค-เนียร์เซนเบอร์ ค่า $k = 1$
ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟีเจอร์แบบคำเดี่ยว ของชุดเอกสาร Reuters-top10

ชุด ที่	k									ค่า สูงสุด
	1	2	3	4	5	6	7	8	9	
1	65.60	71.20	66.40	72.00	69.60	71.20	71.20	73.60	73.60	73.60
2	73.02	81.75	76.19	80.95	78.57	80.95	80.95	80.16	77.78	81.75
3	71.65	73.23	72.44	74.80	75.59	77.17	76.38	76.38	76.38	77.17
4	69.05	70.64	73.81	72.22	72.22	73.02	70.64	72.22	73.02	73.81
ค่าเฉลี่ย	69.83	74.20	72.21	74.99	74.00	75.58	74.79	75.59	75.19	76.58

ตารางที่ 4.48 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เนียบเรสเนเบอร์
ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์แบบคำเดียว ของชุดเอกสาร
Reuters-top10

ชุด ที่	k									ค่า สูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.6480	0.6975	0.6513	0.7021	0.6783	0.6896	0.6873	0.7162	0.7134	0.7162
2	0.7262	0.8062	0.7525	0.7965	0.7742	0.7973	0.7927	0.7869	0.7571	0.8062
3	0.7031	0.7049	0.7001	0.7284	0.7363	0.7476	0.7432	0.7432	0.7442	0.7476
4	0.6742	0.6852	0.7250	0.7137	0.6959	0.7216	0.6907	0.7101	0.7137	0.7250
ค่าเฉลี่ย	0.6879	0.7234	0.7072	0.7352	0.7212	0.7390	0.7285	0.7391	0.7321	0.7488

4.7.4 ผลการทดลองของชุดเอกสาร Reuters-top10เมื่อใช้พีเจอร์แบบคำเดียว และ ไบ-แกรม

สร้างพีเจอร์ที่เป็นคำเดียว และไบ-แกรม โดยสร้างไบ-แกรมจากวลีที่สกัดได้ แล้วทำการเลือกพีเจอร์ที่มีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 1 เอกสาร และทำการทดลองจำแนกประเภทเอกสารโดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.49 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอร์ดังกล่าว ด้วยวิธี เค-เนียบเรสเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.50 และ 4.51 ตามลำดับ

ตารางที่ 4.49 แสดงจำนวนคำเดียวและไบ-แกรมที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง
ของชุดเอกสาร Reuters-top10

ชุดที่	จำนวนพีเจอร์คำเดียวและไบ-แกรม	
	ทั้งหมด	พีเจอร์ที่ถูกเลือก
1	11,043	3,015
2	11,123	2,976
3	10,966	3,098
4	10,975	3,150
ค่าเฉลี่ย	11,027	3,105

ตารางที่ 4.50 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิธีเค-เนียบเรสเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้พีเจอร์แบบคำเดี่ยวและไบ-แกรม ของชุดเอกสาร Reuters-top10

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	68.00	71.20	71.20	73.60	72.00	71.20	72.00	74.40	72.00	74.40
2	76.19	81.75	81.75	81.75	80.95	82.54	83.33	83.33	79.37	83.33
3	74.02	74.80	77.95	76.38	77.95	77.95	78.74	80.32	77.95	80.32
4	72.22	73.81	74.60	75.40	72.22	73.81	72.22	76.98	75.40	76.98
ค่าเฉลี่ย	72.61	75.39	76.38	76.78	75.78	76.38	76.57	78.76	76.18	78.76

ตารางที่ 4.51 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิธีเค-เนียบเรสเนเบอร์ ค่า $k = 1$ ถึง 9 ในแต่ละชุดการเมื่อใช้พีเจอร์แบบคำเดี่ยวและไบ-แกรม ของชุดเอกสาร Reuters-top10

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.6720	0.6949	0.6928	0.7130	0.6932	0.6895	0.6934	0.7246	0.6906	0.7246
2	0.7608	0.8070	0.8074	0.8086	0.7947	0.8111	0.8177	0.8156	0.7738	0.8177
3	0.7264	0.7254	0.7609	0.7371	0.7543	0.7547	0.7684	0.7827	0.7571	0.7827
4	0.7064	0.7280	0.7196	0.7362	0.6941	0.7244	0.7018	0.7526	0.7301	0.7526
ค่าเฉลี่ย	0.7164	0.7388	0.7451	0.7487	0.7341	0.7449	0.7453	0.7689	0.7379	0.7694

4.7.5 ผลการทดลองของชุดเอกสาร Reuters-top10 เมื่อใช้พีเจอร์คำเดี่ยวและวลี

สร้างพีเจอร์จากวลีที่สกัดได้และคำเดี่ยว แล้วทำการเลือกพีเจอร์ที่มีค่าคะแนน (Odds Ratio) มากกว่า 0 และปรากฏในชุดเอกสารทั้งหมดมากกว่า 1 เอกสาร และทำการทดลองจำแนกประเภทเอกสาร โดยในแต่ละชุดการทดลองได้จำนวนพีเจอร์จากเอกสารชุดเรียนรู้ ดังแสดงในตารางที่ 4.52 จากนั้นทำการจำแนกประเภทเอกสารชุดทดสอบ โดยการแทนเอกสารด้วยพีเจอร์ดังกล่าว ด้วยวิธีเค-เนียบเรสเนเบอร์ (k-Nearest Neighbor) ค่า k เท่ากับ 1 ถึง 9 แล้วทำการวัดอัตราความถูกต้องและค่าความถูกต้อง F-measure ดังแสดงในตารางที่ 4.53 และ 4.54 ตามลำดับ

ตารางที่ 4.52 แสดงจำนวนฟิเจอร์คำเดียวและวลี ที่ได้จากชุดเอกสารเรียนรู้ ในแต่ละชุดการทดลอง
ของชุดเอกสาร Reuters-top10

ชุดที่	จำนวนฟิเจอร์คำเดียวและวลี	
	ทั้งหมด	ฟิเจอร์ที่ถูกเลือก
1	9,327	1,882
2	9,424	1,874
3	9,304	1,868
4	9,318	1,965
ค่าเฉลี่ย	9,343	1,897

ตารางที่ 4.53 แสดงอัตราความถูกต้องในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์ ค่า $k = 1$
ถึง 9 ในแต่ละชุดการทดลองเมื่อใช้ฟิเจอร์คำเดียวและวลี ของชุดเอกสาร Reuters-
top10

ชุด ที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	64.80	70.40	68.80	69.60	68.00	69.60	72.00	72.80	72.00	72.80
2	76.98	79.37	76.98	79.37	78.57	77.78	81.75	78.57	79.37	81.75
3	70.87	74.80	72.44	74.02	76.38	77.17	77.17	78.74	77.95	78.74
4	70.64	70.64	72.22	73.81	73.81	73.81	71.43	72.22	73.81	73.81
ค่าเฉลี่ย	70.82	73.80	72.61	74.20	74.19	74.59	75.59	75.58	75.78	76.77

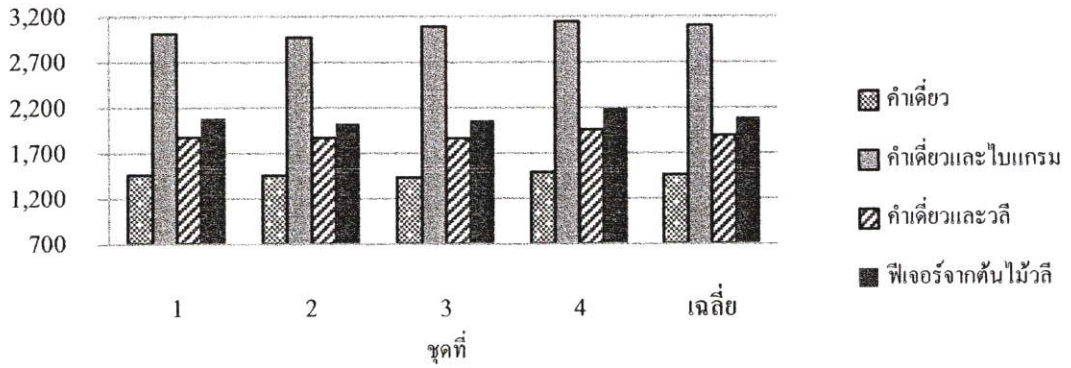
ตารางที่ 4.54 แสดงค่าความถูกต้อง F-measure ในการจำแนกประเภทเอกสารวิเค-เนียร์สเนเบอร์
ค่า $k = 1$ ถึง 9 เมื่อใช้ฟิเจอร์คำเดียวและวลี ของชุดเอกสาร Reuters-top10

ชุดที่	k									ค่าสูงสุด
	1	2	3	4	5	6	7	8	9	
1	0.6393	0.6874	0.6770	0.6837	0.6633	0.6752	0.6951	0.7074	0.7049	0.7074
2	0.7639	0.7814	0.7576	0.7808	0.7742	0.7646	0.8068	0.7743	0.7786	0.8068
3	0.6875	0.7254	0.7024	0.7165	0.7452	0.7486	0.7544	0.7697	0.7609	0.7697
4	0.6966	0.6931	0.7020	0.7253	0.7175	0.7291	0.7065	0.7149	0.7253	0.7291
ค่าเฉลี่ย	0.6968	0.7218	0.7098	0.7266	0.7250	0.7294	0.7407	0.7416	0.7424	0.7533

4.7.6 เปรียบเทียบผลการทดลองเมื่อใช้พีเจอร์แบบต่าง ของชุดเอกสาร Reuters-top10

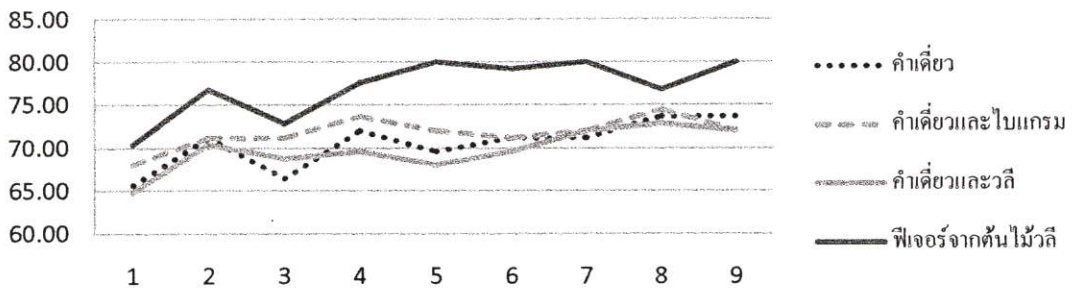
รูปที่ 4.42 ถึง 4.54 แสดงการเปรียบเทียบผลการทดลองเมื่อใช้พีเจอร์แต่ละแบบ โดยเปรียบเทียบทั้งในด้านจำนวนพีเจอร์ อัตราความถูกต้อง และค่า F-measure

จำนวนพีเจอร์ที่ใช้ของชุดเอกสาร Reuters-top10



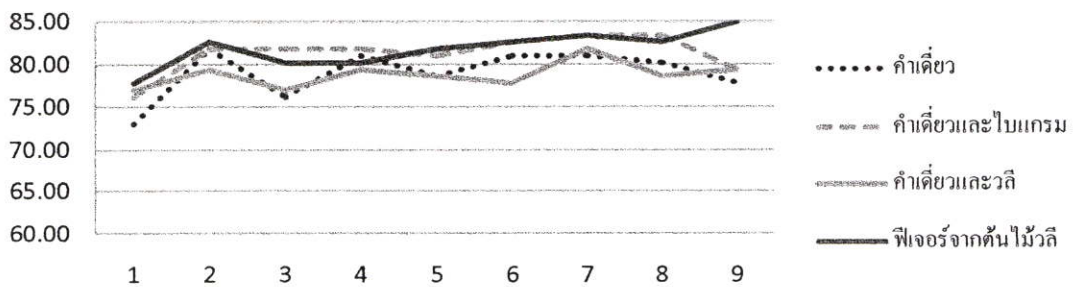
รูปที่ 4.42 แผนภูมิแท่งเปรียบเทียบจำนวนพีเจอร์ที่ใช้ในแต่ละแบบของชุดเอกสาร Reuters-top10

อัตราความถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 1



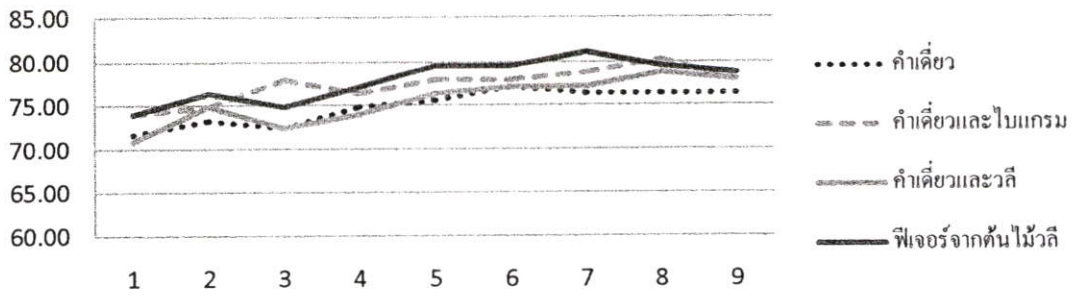
รูปที่ 4.43 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 1

อัตราความถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 2



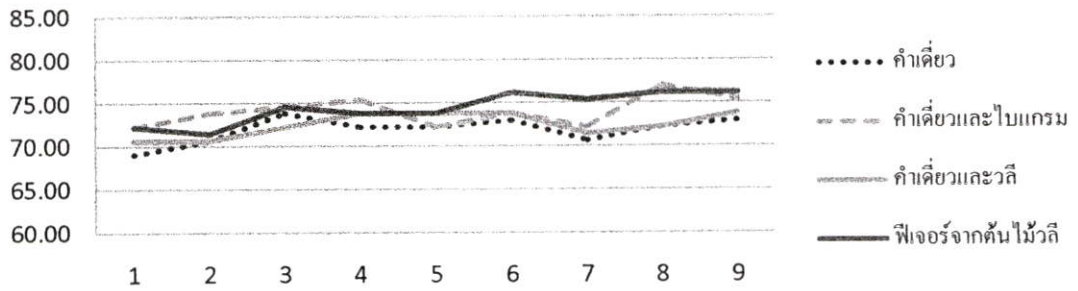
รูปที่ 4.44 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 2

อัตราความถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 3



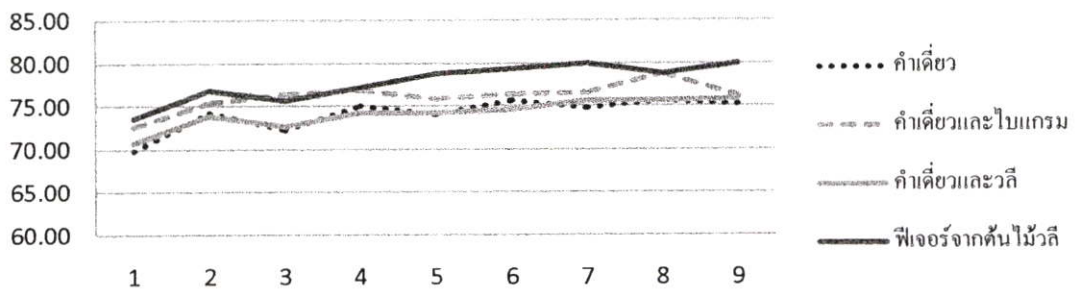
รูปที่ 4.45 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 3

อัตราความถูกต้องเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 4



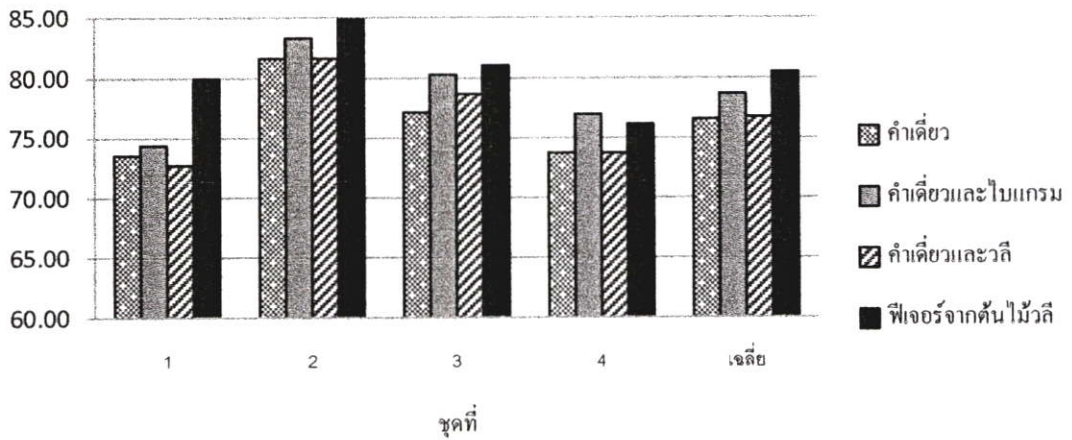
รูปที่ 4.46 กราฟเปรียบเทียบอัตราความถูกต้องในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 4

อัตราความถูกต้องเฉลี่ยเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร Reuters-top10



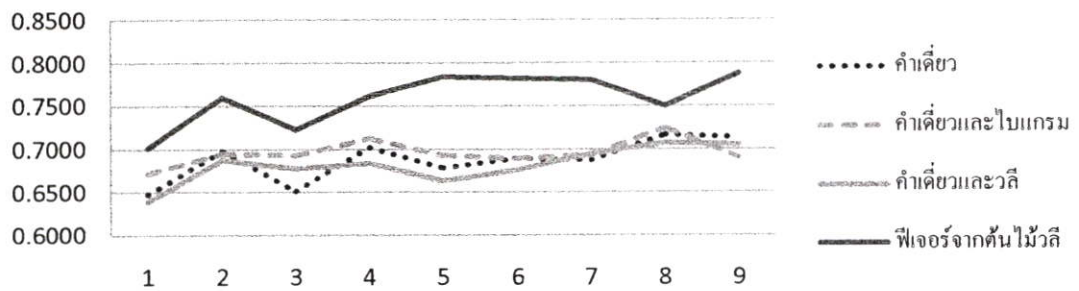
รูปที่ 4.47 กราฟเปรียบเทียบอัตราความถูกต้องเฉลี่ยในแต่ละแบบของชุดเอกสาร Reuters-top10

อัตราความถูกต้องสูงสุดของชุดเอกสาร Reuters-top10



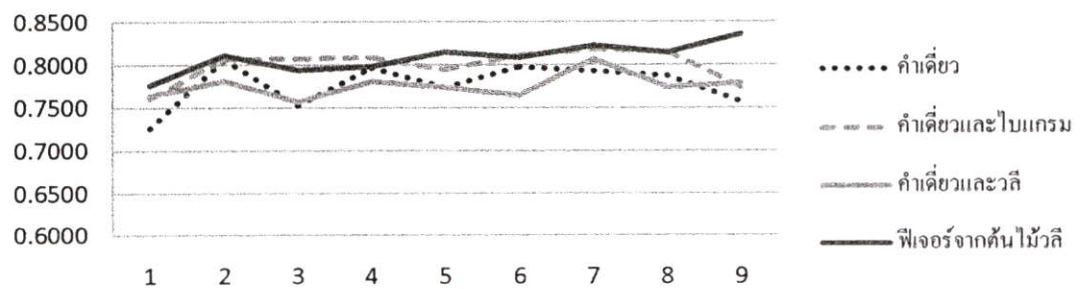
รูปที่ 4.48 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุดในแต่ละแบบของชุดเอกสาร Reuters-top10

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 1



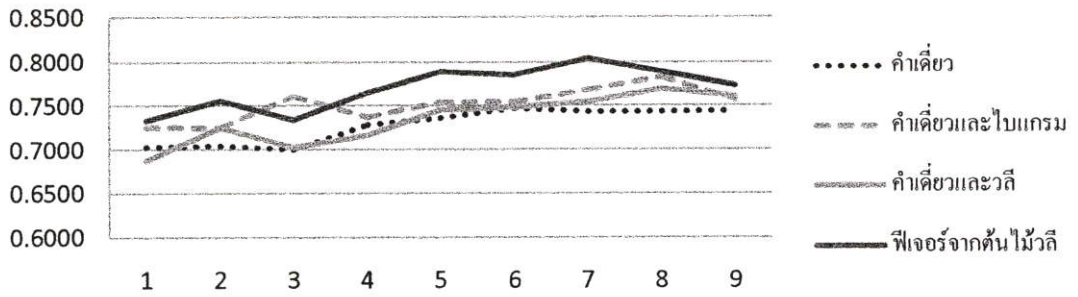
รูปที่ 4.49 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 1

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 2



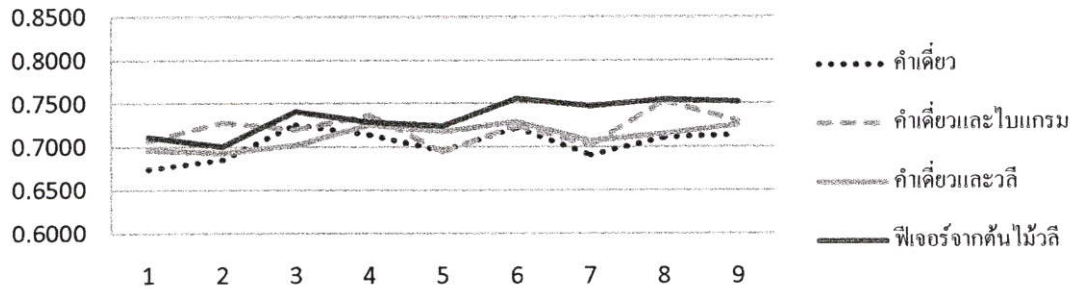
รูปที่ 4.50 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 2

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 3



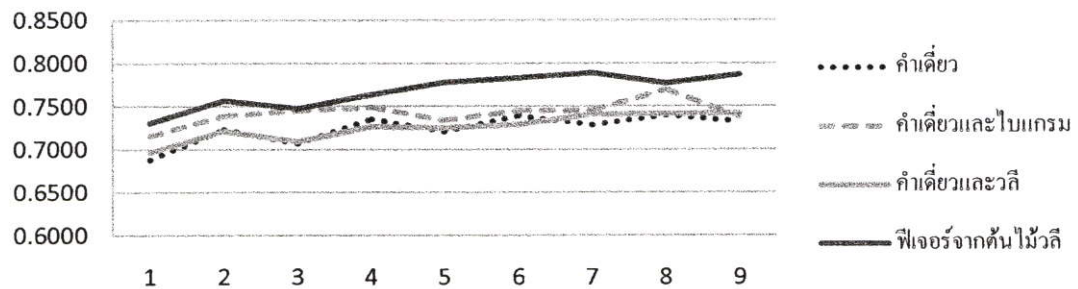
รูปที่ 4.51 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 3

ค่า F-measure เมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 4



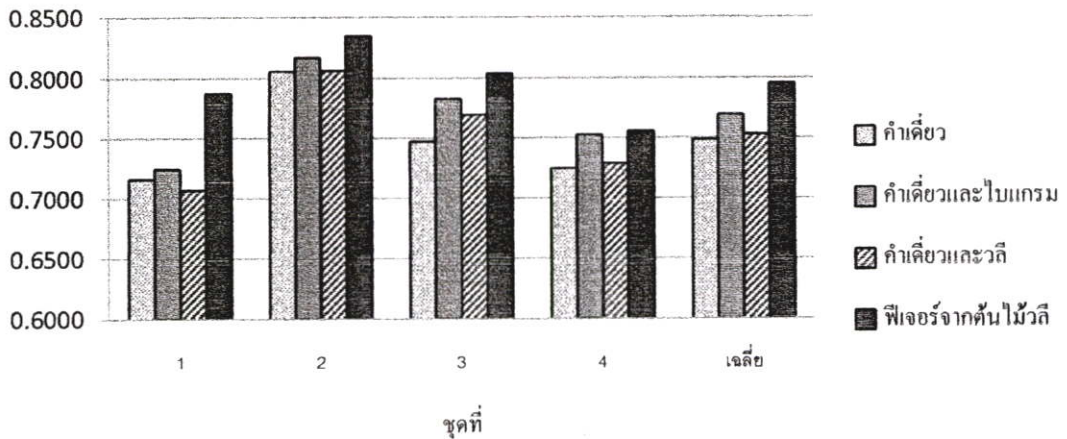
รูปที่ 4.52 กราฟเปรียบเทียบ ค่า F-measure ในแต่ละแบบของชุดเอกสาร Reuters-top10 ชุดที่ 4

ค่า F-measure เฉลี่ยเมื่อใช้ค่า k ต่างๆ ของชุดเอกสาร Reuters-top10



รูปที่ 4.53 กราฟเปรียบเทียบ ค่า F-measure เฉลี่ยในแต่ละแบบของชุดเอกสาร Reuters-top10

ค่า F-measure สูงสุดของชุดเอกสาร Reuters-top10



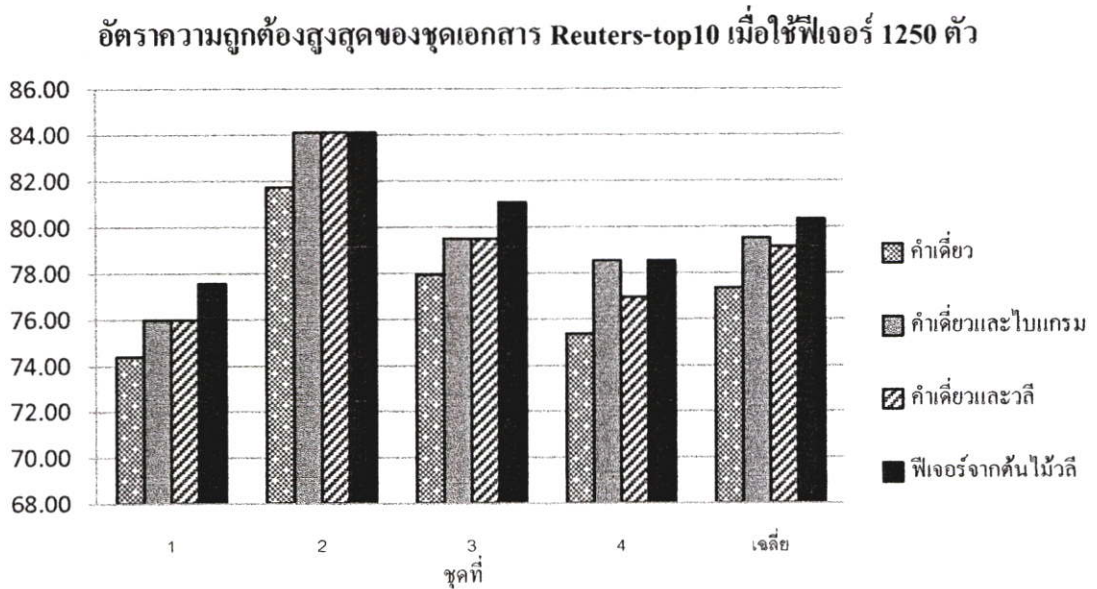
รูปที่ 4.54 แผนภูมิแท่งเปรียบเทียบ ค่า F-measure สูงสุดในแต่ละแบบของชุดเอกสาร Reuters-top10

จากรูปที่ 4.42 จะเห็นได้ว่าฟีเจอร์ที่ได้จากต้นไม้วลีมีจำนวนมากกว่าฟีเจอร์แบบอื่นๆ และเมื่อพิจารณาถึงความถูกต้องในการจำแนกประเภทเอกสาร จากรูปที่ 4.48 และ 4.54 ในทุกชุดการทดลองยกเว้นชุดที่ 4 นั้น ฟีเจอร์จากต้นไม้วลีให้อัตราความถูกต้องและค่า F-measure โดยเฉลี่ยสูงกว่าฟีเจอร์แบบอื่นๆ ซึ่งเป็นไปตามเป้าหมายของเราที่ต้องการได้ฟีเจอร์ที่มีประสิทธิภาพ โดยใช้ลิร่วมกับคำเดียว ที่ให้ความถูกต้องในการจำแนกประเภทเอกสารสูงแต่มีจำนวนน้อย

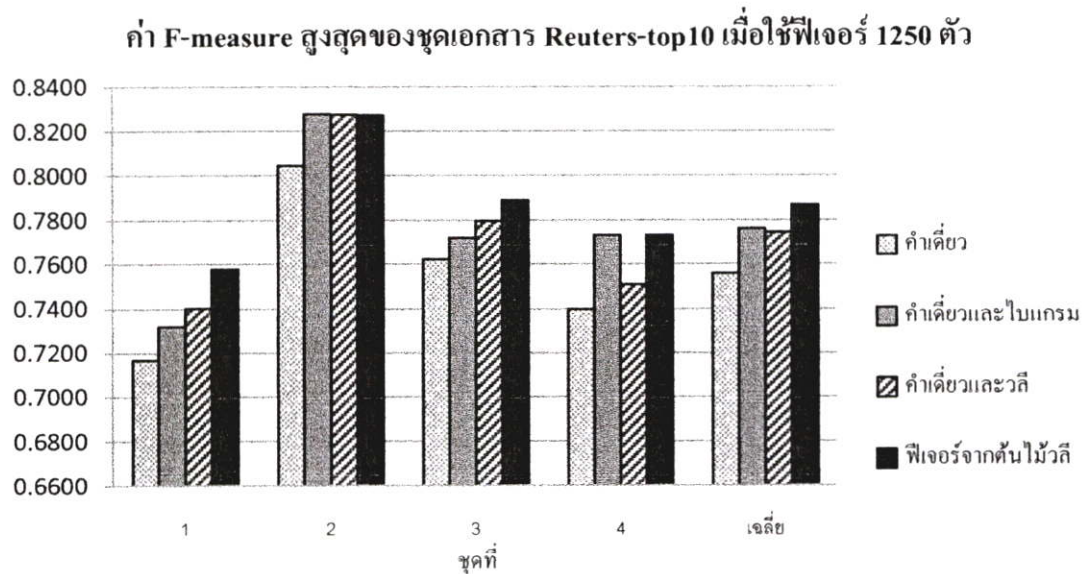
4.7.7 เปรียบเทียบผลการทดลองเมื่อใช้ฟีเจอร์แบบต่าง แบบจำกัดจำนวน ของชุดเอกสาร Reuters-top10

จากผลการทดลองข้างต้น จะเห็นได้ว่าฟีเจอร์จากต้นไม้วลีนั้นทำให้จำแนกประเภทเอกสารได้ความถูกต้องสูงกว่าแบบอื่นๆ แต่ว่ามีจำนวนมากกว่าฟีเจอร์แบบอื่นๆ ยกเว้นฟีเจอร์แบบคำเดียวและไบแกรม ดังนั้น จึงทำการทดลองเอกสารตามแบบข้างต้น แต่เมื่อได้ฟีเจอร์มาแล้วจะทำการจำกัดจำนวนที่จะใช้ฟีเจอร์แต่ละแบบ โดยนำฟีเจอร์มาเรียงตามค่า Odds-Ratio แล้วเลือกฟีเจอร์ที่ให้ค่าคะแนนสูงในลำดับต้นๆ ตามจำนวนที่ต้องการมาใช้ โดยรูปที่ 4.55 ถึงรูปที่ 4.62 แสดงค่าความถูกต้องสูงสุด และค่า F-measure สูงสุดเมื่อใช้จำนวนฟีเจอร์เท่ากับ 1250, 1500, 1750 และ 2000 ตามลำดับ โดยในการเลือกฟีเจอร์ยังคงใช้เกณฑ์ความถี่เอกสารต่ำสุด คือ 1 และค่า Odds-Ratio มากกว่า 0 ผลการเปรียบเทียบเมื่อใช้จำนวนฟีเจอร์ต่างๆ กัน ในแต่ละชุดการทดลอง แสดงดังกราฟในรูปที่ 4.63 ถึง 4.70

รูปที่ 4.71 และ 4.72 แสดงแผนภูมิแท่งเปรียบเทียบค่าอัตราความถูกต้องสูงสุดเฉลี่ย และค่า F-measure สูงสุดเฉลี่ยจากทุกชุดการทดลองเมื่อใช้ฟีเจอร์ต่างๆ ตามลำดับ

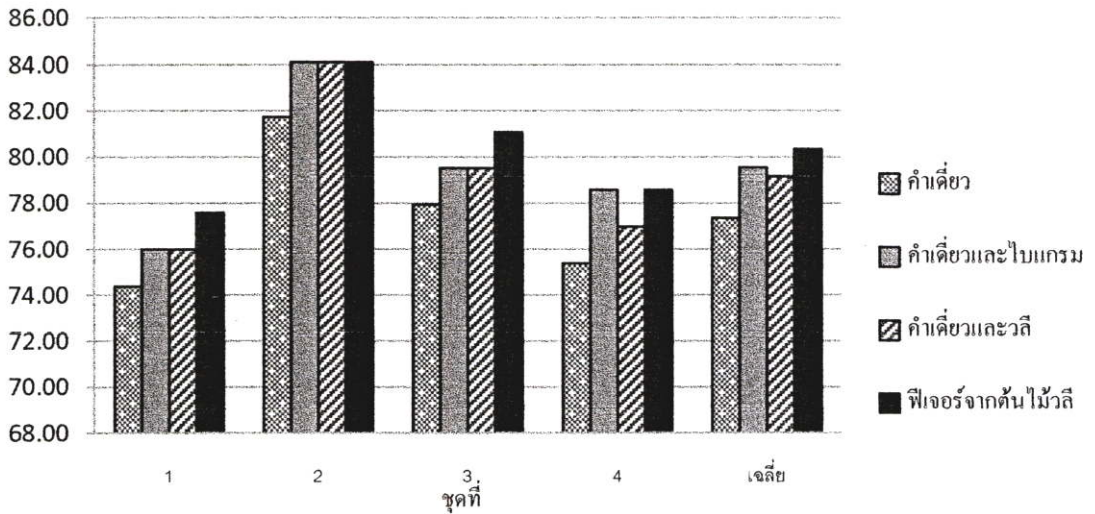


รูปที่ 4.55 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้ฟีเจอร์แต่ละแบบ จำนวน 1250 ตัว ของชุดเอกสาร Reuters-top10



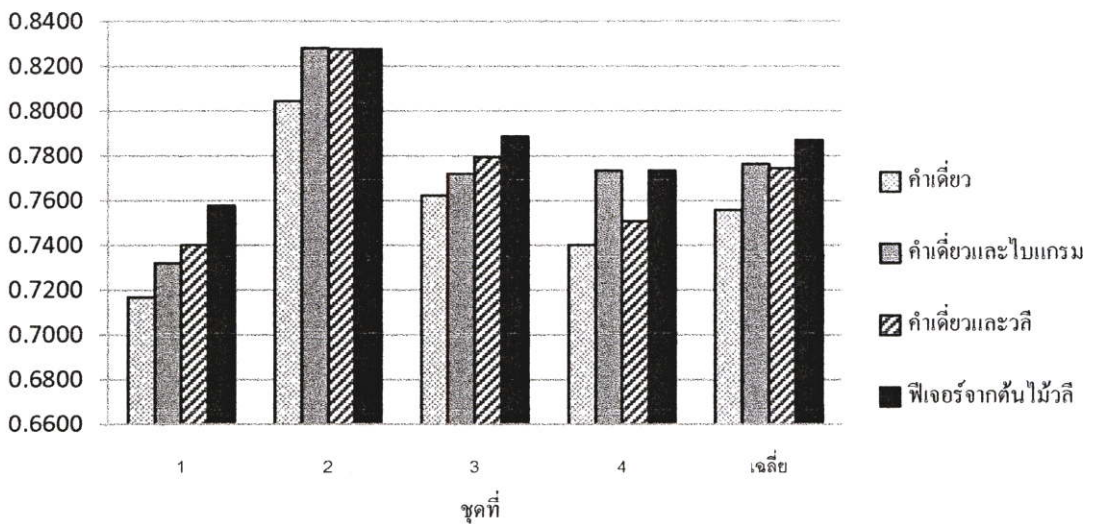
รูปที่ 4.56 แผนภูมิแท่งเปรียบเทียบค่า F-measure สูงสุด เมื่อใช้ฟีเจอร์แต่ละแบบ จำนวน 1250 ตัว ของชุดเอกสาร Reuters-top10

อัตราความถูกต้องสูงสุดของชุดเอกสาร Reuters-top10 เมื่อใช้ฟิเจอร์ 1500 ตัว



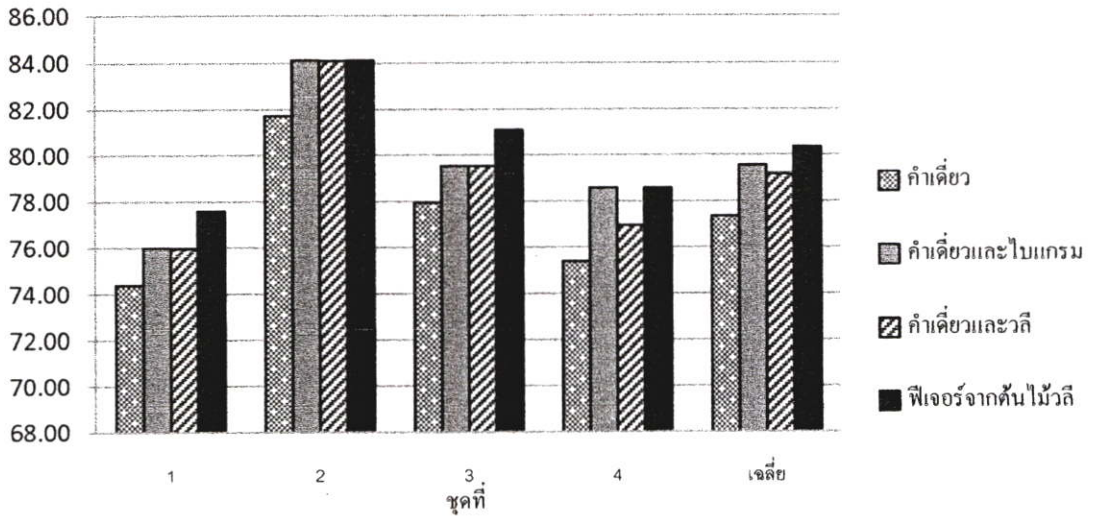
รูปที่ 4.57 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 1500 ตัว ของชุดเอกสาร Reuters-top10

ค่า F-measure สูงสุดของชุดเอกสาร Reuters-top10 เมื่อใช้ฟิเจอร์ 1500 ตัว



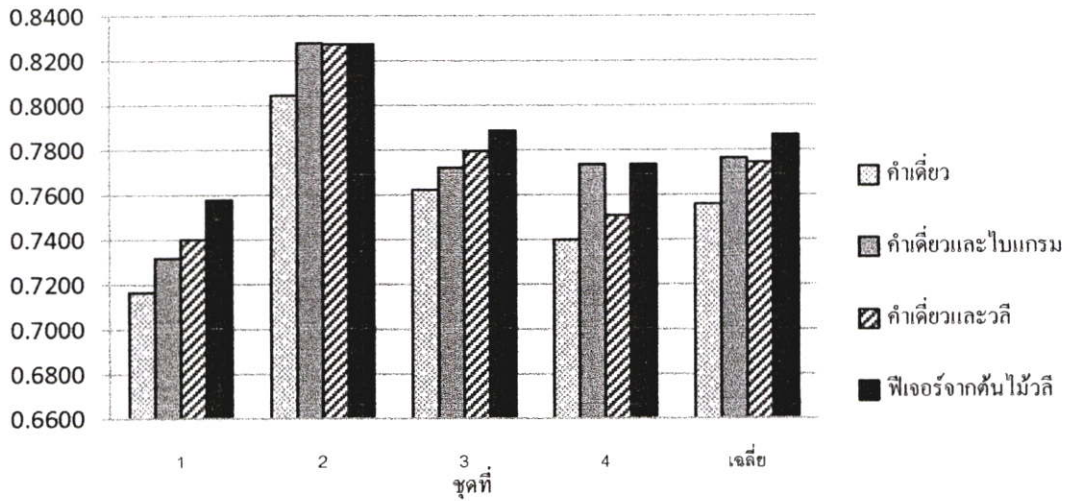
รูปที่ 4.58 แผนภูมิแท่งเปรียบเทียบ ค่า F-measure สูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 1500 ตัว ของชุดเอกสาร Reuters-top10

อัตราความถูกต้องสูงสุดของชุดเอกสาร Reuters-top10 เมื่อใช้ฟิวเจอร์ 1750 ตัว



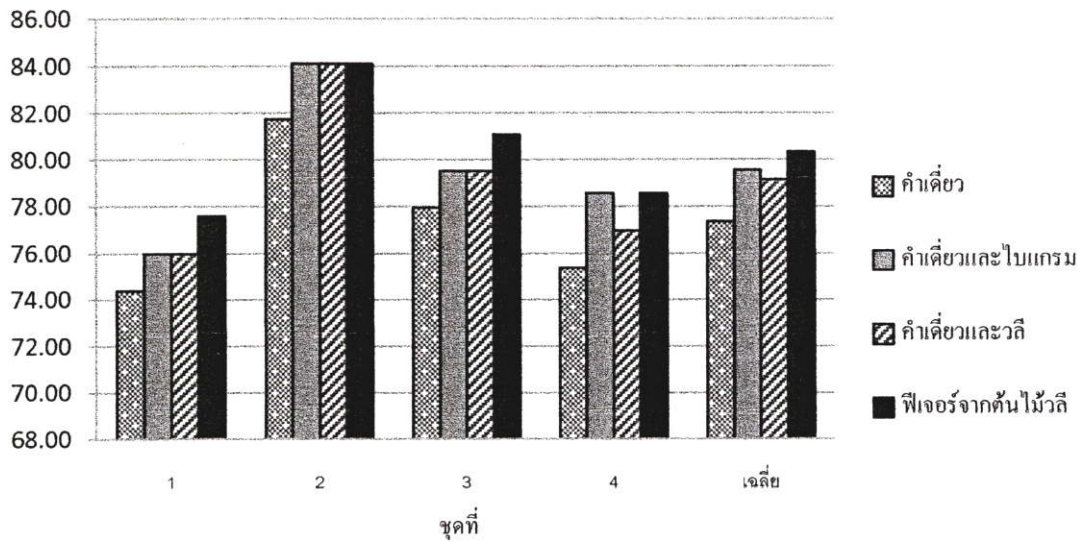
รูปที่ 4.59 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้ฟิวเจอร์แต่ละแบบ จำนวน 1750 ตัว ของชุดเอกสาร Reuters-top10

ค่า F-measure สูงสุดของชุดเอกสาร Reuters-top10 เมื่อใช้ฟิวเจอร์ 1750 ตัว



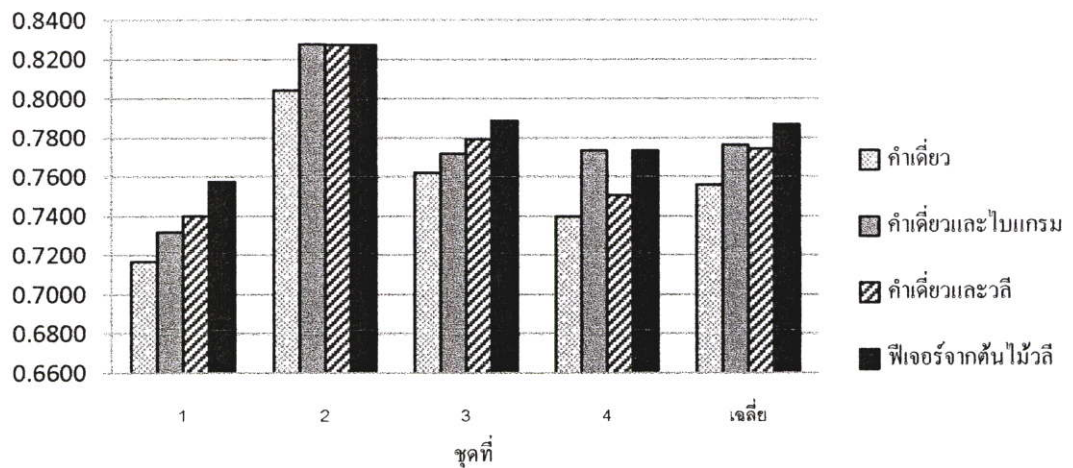
รูปที่ 4.60 แผนภูมิแท่งเปรียบเทียบ ค่า F-measure สูงสุด เมื่อใช้ฟิวเจอร์แต่ละแบบ จำนวน 1750 ตัว ของชุดเอกสาร Reuters-top10

อัตราความถูกต้องสูงสุดของชุดเอกสาร Reuters-top10 เมื่อใช้ฟิเจอร์ 2000 ตัว



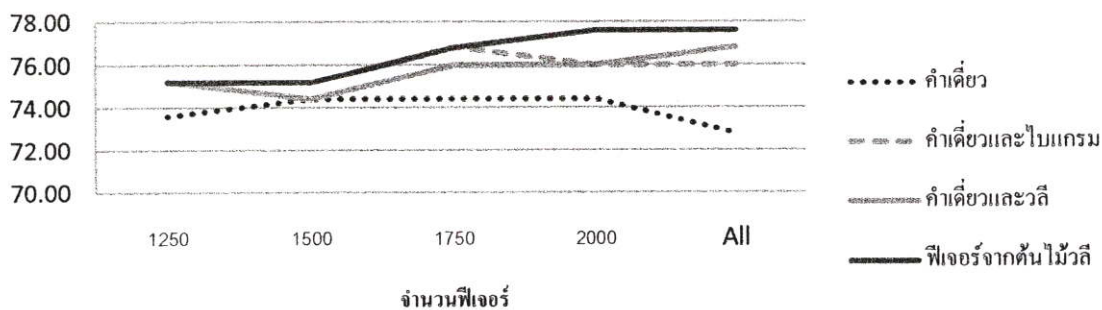
รูปที่ 4.61 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 2000 ตัว ของชุดเอกสาร Reuters-top10

ค่า F-measure สูงสุดของชุดเอกสาร Reuters-top10 เมื่อใช้ฟิเจอร์ 2000 ตัว



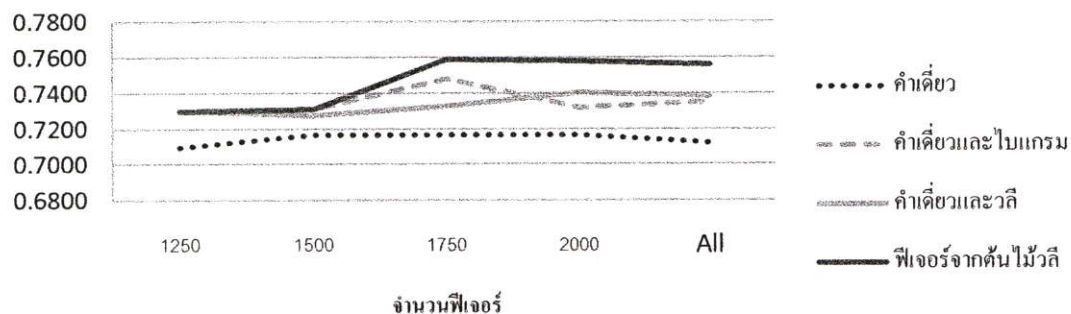
รูปที่ 4.62 แผนภูมิแท่งเปรียบเทียบค่า F-measure สูงสุด เมื่อใช้ฟิเจอร์แต่ละแบบ จำนวน 2000 ตัว ของชุดเอกสาร Reuters-top10

อัตราการถูกต้องเมื่อใช้จ.น.พีเจอร์ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 1



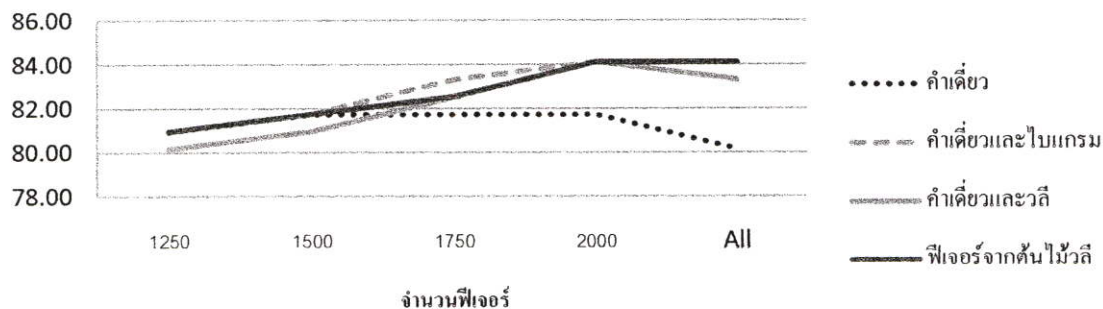
รูปที่ 4.63 กราฟเปรียบเทียบอัตราการถูกต้องสูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 1

ค่า F-measure เมื่อใช้จ.น.พีเจอร์ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 1



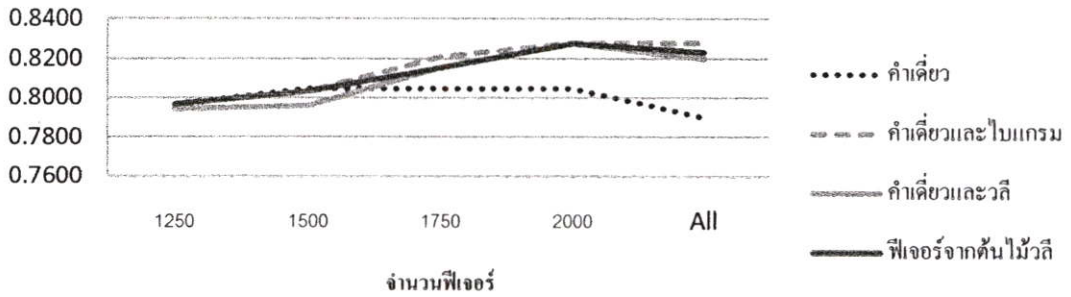
รูปที่ 4.64 กราฟเปรียบเทียบค่า F-measure สูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 1

อัตราการถูกต้องเมื่อใช้จ.น.พีเจอร์ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 2



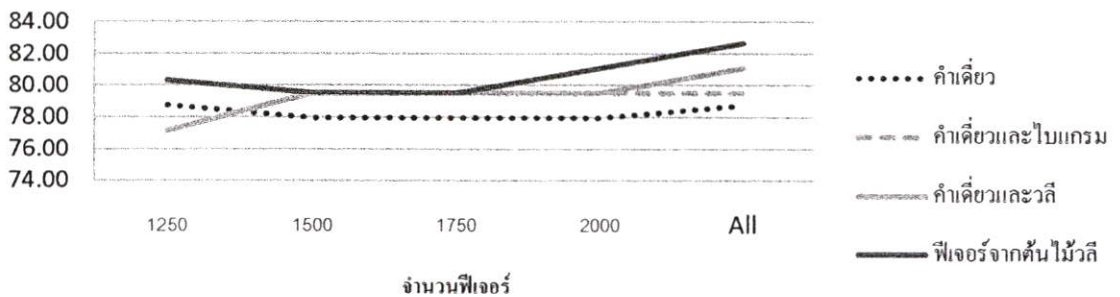
รูปที่ 4.65 กราฟเปรียบเทียบอัตราการถูกต้องสูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 2

ค่า F-measure เมื่อใช้จ.น.พีเจอร์ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 2



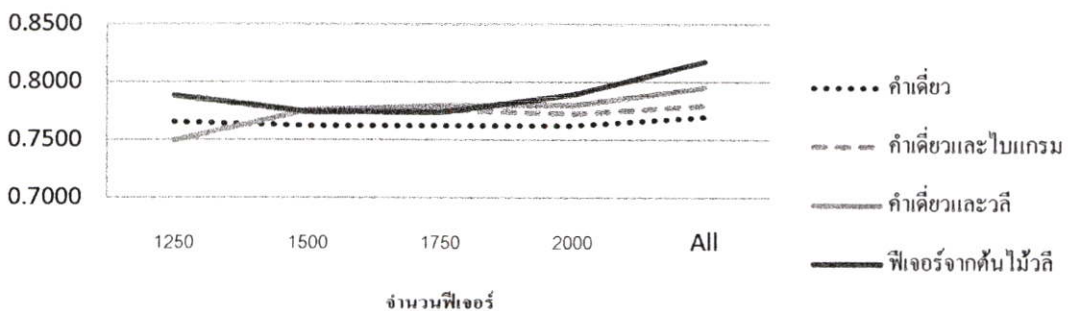
รูปที่ 4.66 กราฟเปรียบเทียบค่า F-measure สูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 2

อัตราการถูกต้องเมื่อใช้จ.น.พีเจอร์ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 3



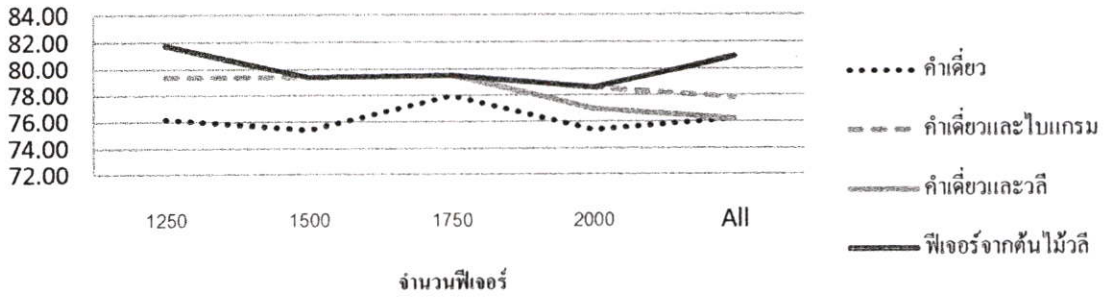
รูปที่ 4.67 กราฟเปรียบเทียบอัตราการถูกต้องสูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 3

ค่า F-measure เมื่อใช้จ.น.พีเจอร์ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 3



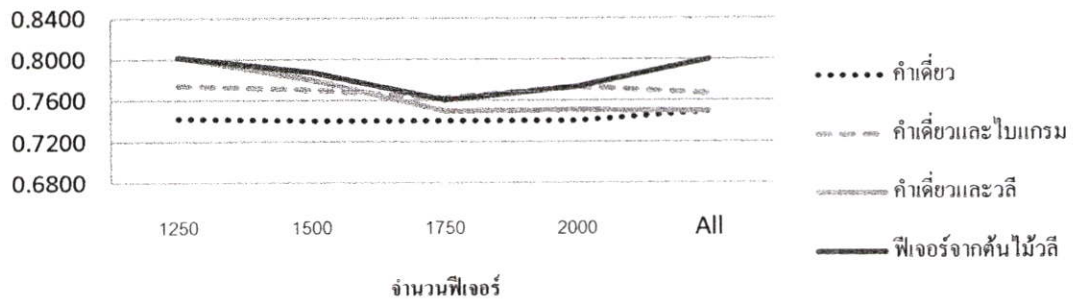
รูปที่ 4.68 กราฟเปรียบเทียบค่า F-measure สูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 3

อัตราความถูกต้องเมื่อใช้ช.น.พีเจอร์ท่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 4



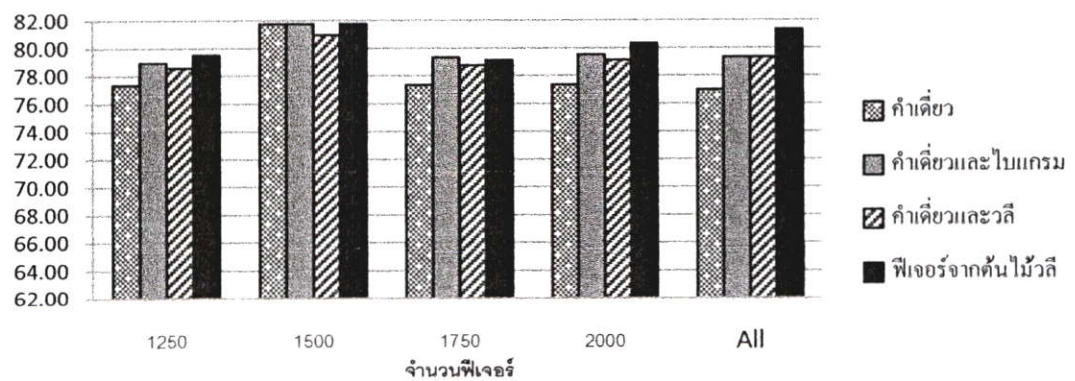
รูปที่ 4.69 กราฟเปรียบเทียบอัตราความถูกต้องสูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 4

ค่า F-measure เมื่อใช้ช.น.พีเจอร์ต่างๆ ของชุดเอกสาร Reuters-top10 ชุดที่ 4



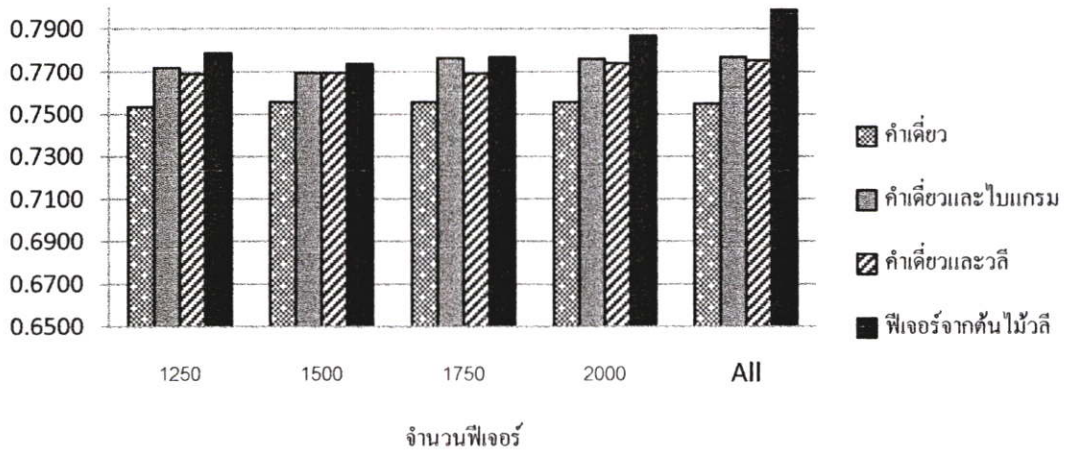
รูปที่ 4.70 กราฟเปรียบเทียบค่า F-measure สูงสุด เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10 ชุดที่ 4

อัตราความถูกต้องเฉลี่ยเมื่อใช้ช.น.พีเจอร์ต่างๆ ของชุดเอกสาร Reuters-top10



รูปที่ 4.71 แผนภูมิแท่งเปรียบเทียบค่าอัตราความถูกต้องเฉลี่ยเมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10

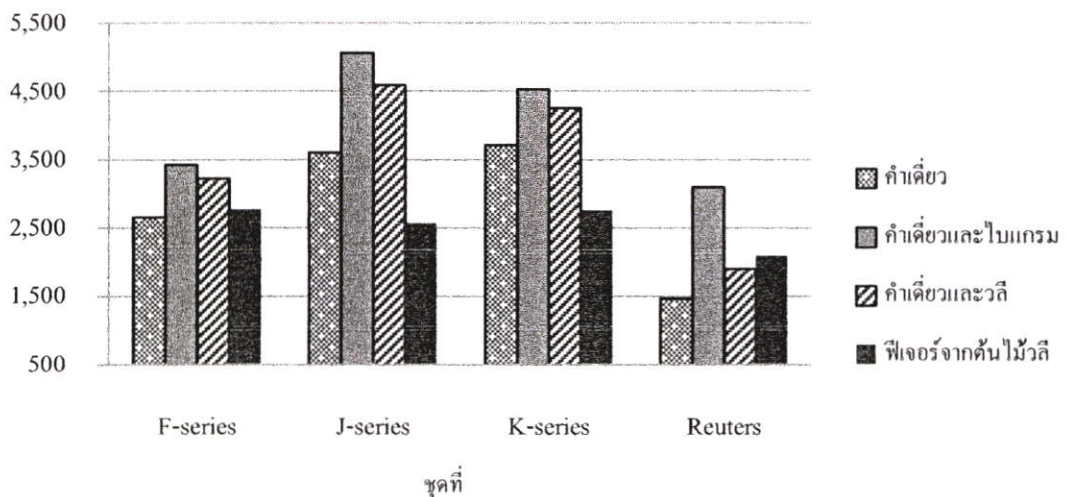
ค่า F-measure เฉลี่ยเมื่อใช้จ.น.พีเจอร์ต่างๆ ของชุดเอกสาร Reuters-top10



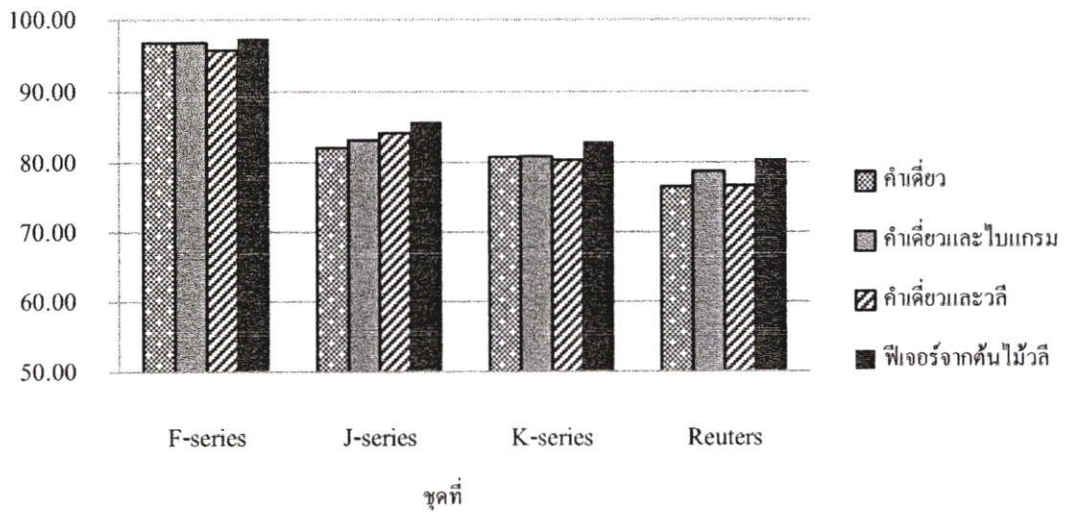
รูปที่ 4.72 แผนภูมิแท่งเปรียบเทียบค่า F-measure เฉลี่ย เมื่อใช้พีเจอร์ในแต่ละแบบ โดยมีการจำกัดจำนวนของพีเจอร์ ของชุดเอกสาร Reuters-top10

4.8 การเปรียบเทียบผลการทดลองทุกชุดข้อมูล

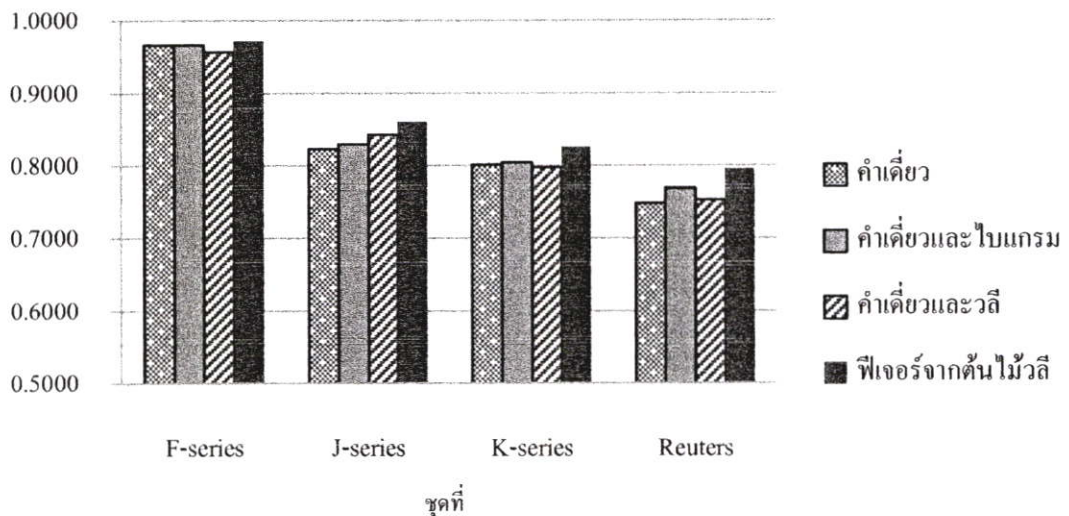
จากผลการทดลองในชุดข้อมูล PDDP ทั้ง F-series, J-series และ K-series การลดพีเจอร์โดยใช้ต้นไม้วลีนั้นสามารถลดจำนวนของพีเจอร์ และเพิ่มความถูกต้องได้ แต่ในข้อมูล Reuters-top10 จำนวนพีเจอร์ยังคงมากกว่าพีเจอร์แบบอื่นๆ แต่สามารถเพิ่มความถูกต้องในการจำแนกประเภทเอกสารได้ ดังแสดงในรูปที่ 4.55, 4.56 และ 4.57 ซึ่งแสดงผลการเปรียบเทียบจำนวนพีเจอร์ อัตราความถูกต้อง และค่า F-measure เมื่อใช้พีเจอร์แบบต่างๆ ของทุกชุดข้อมูลตามลำดับ



รูปที่ 4.73 แผนภูมิแท่งเปรียบเทียบจำนวนพีเจอร์แต่ละแบบของทุกชุดข้อมูล



รูปที่ 4.74 แผนภูมิแท่งเปรียบเทียบอัตราความถูกต้องเมื่อใช้ฟิเจอร์แต่ละแบบของทุกชุดข้อมูล



รูปที่ 4.75 แผนภูมิแท่งเปรียบเทียบค่า F-measure เมื่อใช้ฟิเจอร์แต่ละแบบของทุกชุดข้อมูล

พิจารณาจำนวนฟิเจอร์จากรูปที่ 4.73 จะเห็นได้ว่าในชุดข้อมูล J-series และ K-series ฟิเจอร์ที่ได้จากต้นไม้วลีมีปริมาณน้อยกว่าฟิเจอร์แบบคำเดียว ทั้งนี้เนื่องจากเราใช้เกณฑ์ความถี่เอกสารในการเลือกวลีที่สกัดได้ ก่อนนำมาสร้างต้นไม้วลี ทำให้บางวลีไม่ได้ถูกเลือก และทำให้ปริมาณคำเดียวลดลง ทำให้มีปริมาณน้อยกว่าฟิเจอร์แบบคำเดียว

จากการทดลองทั้ง 4 ชุดข้อมูล จะเห็นได้ว่าฟิเจอร์จากต้นไม้วลีมีประสิทธิภาพ คือ ให้ความถูกต้องในการจำแนกประเภทเอกสารมากกว่าฟิเจอร์แบบอื่นๆ และโดยส่วนใหญ่แล้วได้จำนวนที่น้อยกว่าฟิเจอร์แบบอื่นๆ ด้วย

ในบทที่ 5 จะกล่าวถึงผลสรุปของงานวิจัยและข้อเสนอแนะ ซึ่งจะเป็นประโยชน์ในการพัฒนาการทำวิจัยต่อไป

บทที่ 5

สรุปผลการวิจัยและข้อเสนอแนะ

5.1 สรุปผลการวิจัย

การจำแนกประเภทเอกสารนั้น คือ การจัดประเภทเอกสารให้กับเอกสารใหม่ โดยใช้ความรู้จากเอกสารเดิมที่เราทราบประเภทอยู่แล้วมาใช้ การจำแนกประเภทเอกสารให้มีประสิทธิภาพ คือ มีความถูกต้องสูง และใช้เวลาไม่นาน ปัจจัยที่สำคัญก็คือ ฟีเจอร์ที่ใช้แทนเอกสาร ซึ่งจะต้องเป็นตัวแทนเอกสารได้ดี คือ สื่อถึงความหมายของประเภทเอกสารต่างๆ ได้ และจะต้องมีจำนวนที่ไม่มากเกินไป เพื่อลดเวลาในการประมวลผล

งานวิจัยนี้ได้นำเสนอการลดฟีเจอร์สำหรับการจำแนกประเภทเอกสาร โดยใช้ต้นไม้ตัดสินใจ ซึ่งเป็นการใช้โครงสร้างข้อมูลแบบต้นไม้มาช่วยในการสร้างฟีเจอร์จากวลีที่สกัดได้ และช่วยในการเลือกฟีเจอร์โดยดูจากความสัมพันธ์ของฟีเจอร์ภายในต้นไม้ โดยใช้ค่า Odds-Ratio ร่วมด้วย โดยได้ทำการทดลองวัดประสิทธิภาพของฟีเจอร์ที่ได้จากต้นไม้ตัดสินใจกับฟีเจอร์แบบอื่นๆ โดยดูจากจำนวนและความถูกต้องในการจำแนกประเภทเอกสารที่ได้ โดยใช้ชุดเอกสารตัวอย่างในการทดสอบ ซึ่งเป็นชุดเอกสารที่ถูกใช้ในหลายๆ งานวิจัยทางด้าน การจำแนกประเภทเอกสาร

ซึ่งจากผลการทดลองสรุปได้ว่าการใช้ต้นไม้ตัดสินใจในการลดฟีเจอร์ ทำให้ได้ฟีเจอร์ที่มีประสิทธิภาพคือ มีจำนวนน้อยแต่ยังคงความถูกต้องในการจำแนกประเภทเอกสารไว้ได้ และสามารถนำงานวิจัยนี้ไปประยุกต์ใช้ในงานต่างๆ ได้ เช่น การจำแนกประเภทของจดหมายอิเล็กทรอนิกส์โดยแบ่งเป็นประเภทที่เป็นขยะ และ ไม่เป็นขยะ ซึ่งจะนำเนื้อหาของแต่ละจดหมายอิเล็กทรอนิกส์มาพิจารณา หรืออาจนำไปประยุกต์ใช้ในการจำแนกหนังสือในห้องสมุดที่มีการแบ่งประเภทไว้แล้ว โดยอาจจะนำส่วนของชื่อเรื่อง คำนำ และสารบัญ มาใช้ในการพิจารณาจำแนกประเภท แทนการพิจารณาในส่วนเนื้อหาของเนื้อหาทั้งหมดของหนังสือ

5.2 ข้อเสนอแนะ

จากปัญหาที่จำนวนฟีเจอร์ที่ได้จากการเลือกโดยต้นไม้ตัดสินใจมีจำนวนมากกว่าแบบอื่นๆ ในชุดเอกสาร Reuters-top10 จึงควรหาวิธีที่จะช่วยลดความซ้ำซ้อนของฟีเจอร์ในต้นไม้ตัดสินใจให้ได้มากกว่านี้ แต่ยังคงรักษาความถูกต้องในการจำแนกประเภทเอกสารไว้

ในเรื่องของการเลือกฟีเจอร์ในต้นไม้ตัดสินใจในงานวิจัยนี้ใช้เพียงแค่ว่า Odds Ratio ในการให้คะแนนในแต่ละเทอม ควรหาเกณฑ์อื่นมาใช้ด้วย เช่น อาจจะมีการถ่วงน้ำหนักคำที่มีความสำคัญ คำที่อยู่ในส่วนหัวข้อ เป็นต้น หรืออาจใช้ค่าคะแนนอื่นที่ไม่ใช่ Odds Ratio ก็ได้

บรรณานุกรม

- [1] T. Dumis, J. Platt, D. Heckerman, and M. Sahami. "**Inductive learning algorithms and representation for text categorization**". CIKM'98, ACM Info. & KM. Bethesda, US, pp. 148-155, 1998.
- [2] G. Salton, A. Wong, and C. S. Yang, "**A Vector Space Model for Automatic Indexing**," Communications of the ACM, vol. 18, nr. 11, pages 613–620, 1975.
- [3] F. Sebastiani and Franca Debole. "**Supervised term weighting for automated text categorization**". SAC-03 18th, ACM Applied Computing, pp. 784-788, 2003.
- [4] F. Sebastiani. "**Machine learning in automated text categorization**". ACM Computing Surveys, 34(1), pp. 1-47, March 2002.
- [5] G. Gonged, H. Wang, D. Bell, Y. Bi, and K. Gree. "**Using kNN model-based approach for automatic text categorization**". European Commission project ICONS, project no. IST-2001-32429.
- [6] M. Sholom, T. Zhang and J. Fred, **Text Mining: Predictive Methods for Analyzing Unstructured Information**. MA, USA., Springer, 2004.
- [7] M.F. Porter. "**An Algorithm for Suffix Stripping**". Program, 14 no.3, pp. 130-137, 1980.
- [8] C.E. Shannon. "**A Mathematical Theory of Communication**", Bell System Technical Journal, 27, pp. 379–423 & 623–656, July & October, 1948.
- [9] Y. Y. Yao. **Information-theoretic measures for knowledge discovery and dat mining, in Entropy measures, maximum entropy principle and emerging applications**, Karmeshu (ed.), Springer, pp. 115-136, 2003.
- [10] M. F. Caroppreso, S. Matwin, and F. Sebastani. "**A learner-independent evaluation of the usefulness of statistical phrases for automated text categorization**", Text Databases and Document Management: Theory and Practice, Hershey, PA, pp. 78-102, 2001.
- [11] G.Salton, and C.Buckley, "**Term-weighting approaches in automatic text retrieval**". Information Processing & Management 24(5): pp. 513–523, 1988.
- [12] P.C. Mahalanobis, "**On the generalised distance in statistics**", Proceedings of the National Institute of Science of India 12, pp. 49-55, 1936.

- [13] G. Salton, **Automatic Text Processing: the transformation, analysis and retrieval of information by computer**. New York: Addison Wesley Publishing Company. 1989.
- [14] Rish, Irina. "An empirical study of the naive Bayes classifier". IJCAI 2001 Workshop on Empirical Methods in Artificial Intelligence, 2001.
- [15] Baeza-Yates, R. and Ribeiro-Neto, B. **Modern Information Retrieval**. New York : Addison Wesley Publishing Company. 1999.
- [16] **PDDP Data Set**, <http://ftp.cs.umn.edu/dept/users/boley/PDDPdata/>.
- [17] Adam Schenker, Mark Last, Horst Bunke and Abraham Kandelm, "**Classification of Web Documents Using a Graph Model**". Proceeding of the 7th ICDAR 2003, IEEE, 2003.
- [18] D.D. Lewis, "**Reuters-21578 text categorization test collection distribution 1.0.**", 2004. [Online]. Available : <http://www.daviddlewis.com/resources/testcollections/reuters21578/>.
- [19] [Online]. Available : <http://www.recltwo.com/datasets.html>.
- [20] สุรินทร์ นิยมางกูร, **เทคนิคการสุ่มตัวอย่าง**, สำนักพิมพ์มหาวิทยาลัยเกษตรศาสตร์, พิมพ์ครั้งที่ 3, กรุงเทพฯ, 2541.
- [21] J. Klok, S. Driessen, and M. Brunner, "**Some terms are more interchangeable than others**". Cross-Language Evaluation Forum (CLEF) 2001.
- [22] Stopwords List [Online]. Available : http://www.dcs.gla.ac.uk/idom/ir_resources/linguistic_utils/stop_words.
- [23] WordNet [Online]. Available : <http://wordnet.princeton.edu/online/>.
- [24] D.D. Lewis, "**An Evaluation of Phrasal and Clustered Representations on a Text Categorization Task**". Proceeding of SIGIR'92, pp. 37-50, ACM Press, New York, 1992.
- [25] D. Mladeni and M. Grobelink, "**Word sequences as features in text learning**". 17th Electro. & Comp. Science Conf., Slovenia, pp. 145-148, 1998.
- [26] C. M. Tan, Y.F. Wang and C. D. Lee, "**The use of bigrams to enhance text categorization**". Info. Proc. & Management, 38(4): pp. 529-546, 2002.
- [27] R. Bekkerman and J. Allan, "**Using Bigrams in Text Categorization**". CIIR Technical Report IR-408, University of Massachusetts, 2003.

ภาคผนวก

ภาคผนวก ก.

แสดงตัวอย่างชุดเอกสารต่างๆ ที่ใช้ในการทดลอง

ก.1 ชุดเอกสาร PDDP

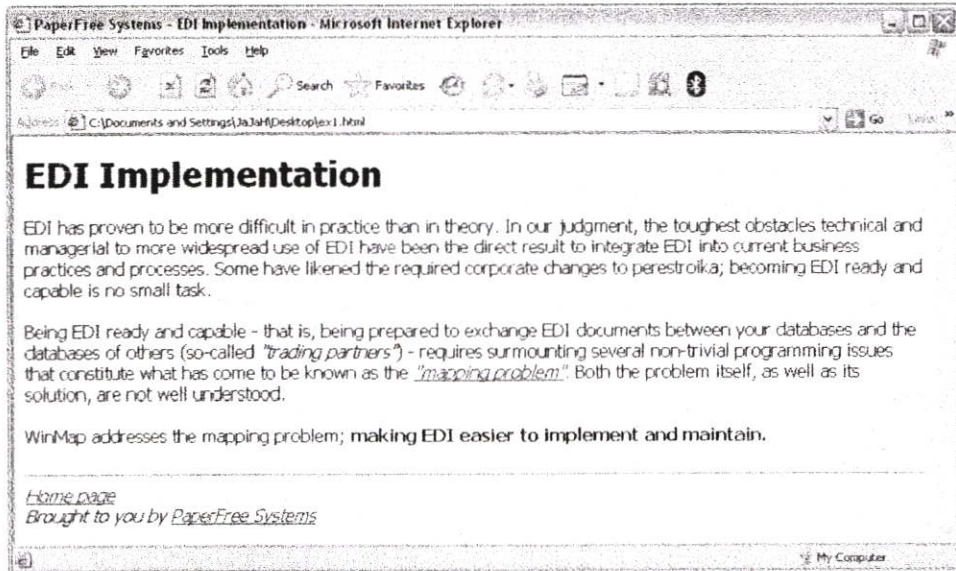
ชุดเอกสาร PDDP [16] เป็นเอกสารเว็บ ซึ่งแบ่งเป็น 3 ชุดย่อย F-series J-series และ K-series โดยแต่ละเอกสารอยู่ในรูปแบบเอกสาร HTML ทั้ง 3 ชุดย่อยมีลักษณะเอกสารเหมือนกัน และมีรายละเอียดและขั้นตอนการเตรียมเอกสารต่างๆ ดังนี้

ก.1.1 ตัวอย่างไฟล์เอกสาร PDDP F-series ในรูปแบบ HTML

รูปแบบของเอกสารก่อนการตัดแท็ก เป็นดังนี้

```
<!-- www.paperfree.com/ediprob.htm -->
<html><head><title>PaperFree Systems - EDI Implementation</title></head>
<body><!-- ediprob.htm --><h1><B>EDI Implementation</B></h1>
<p>EDI has proven to be more difficult in practice than in theory. In our judgment, the toughest obstacles technical and managerial to more widespread use of EDI have been the direct result to integrate EDI into current business practices and processes. Some have likened the required corporate changes to perestroika; becoming EDI ready and capable is no small task.
<p>Being EDI ready and capable - that is, being prepared to exchange EDI documents between your databases and the databases of others (so-called <i>"trading partners"</i>) - requires surmounting several non-trivial programming issues that constitute what has come to be known as the <a href="/edi/challenge.htm"><i>"mapping problem"</i></a>. Both the problem itself, as well as its solution, are not well understood.
<p>WinMap addresses the mapping problem; <b>making EDI easier to implement and maintain.</b> <p><hr>
<i><a href="/edi/index.htm">Home page </A></i>
<br>
<i>Brought to you by <a href="/edi/pfs.htm">PaperFree Systems</A></i>
</BODY>
</html>
```

หน้าตาเว็บเพจของ เอกสารเว็บตัวอย่างนี้ ดังรูป ก.1



รูปที่ ก.1 ตัวอย่างเว็บเพจของเอกสารตัวอย่าง PDDP F-series

ก.1.2 ตัวอย่างไฟล์เอกสาร PDDP F-series ที่ตัดแท็กออกแล้ว

PaperFree Systems - EDI Implementation.

EDI Implementation.

.EDI has proven to be more difficult in practice than in theory. In our judgment, the toughest obstacles technical and managerial to more widespread use of EDI have been the direct result of . to integrate EDI into current business practices and processes. Some have likened the required corporate changes to perestroika; becoming EDI ready and capable is no small task.

.Being EDI ready and capable - that is, being prepared to exchange EDI documents between your databases and the databases of others (so-called .trading partners")- requires surmounting several non-trivial programming issues that constitute what has come to be known as the ."mapping problem".

Both the problem itself, as well as its solution, are not well understood.

.WinMap addresses the mapping problem; .making EDI easier to implement and maintain..

.Brought to you by .

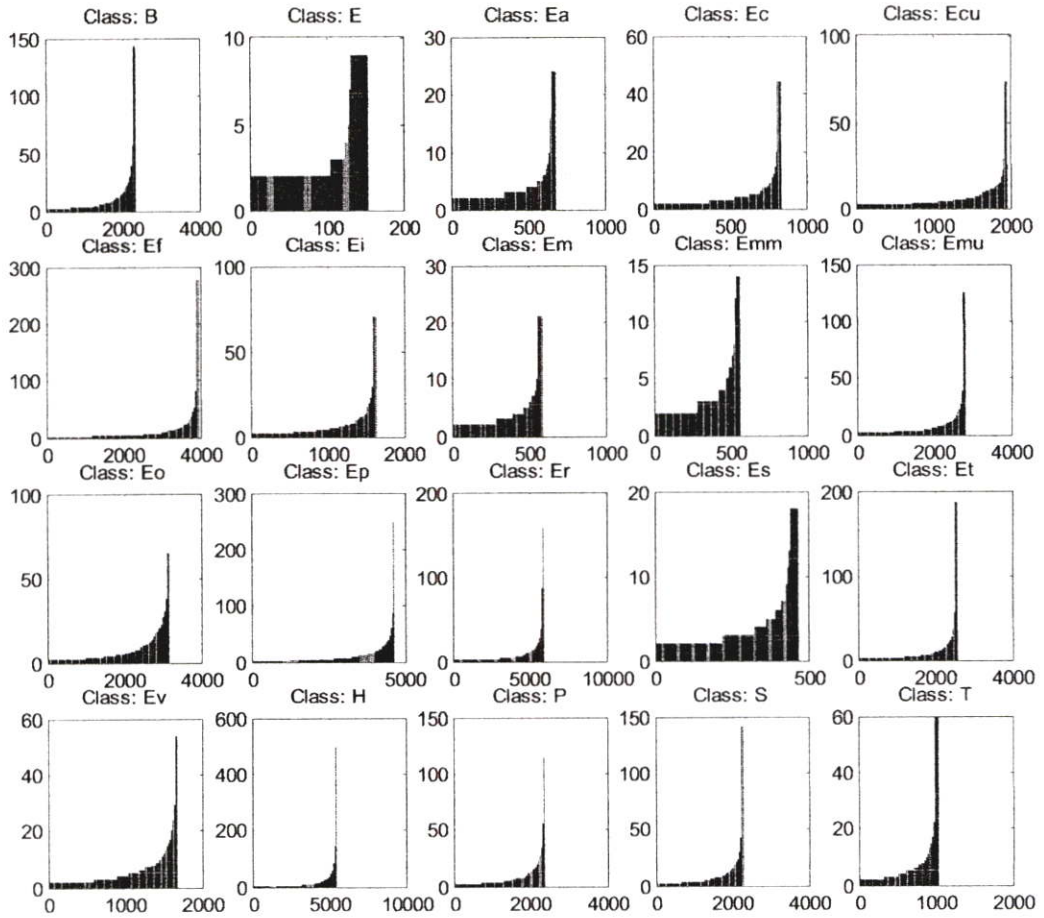
ก.1.3 วิธีที่สกัดได้จากตัวอย่างไฟล์เอกสาร PDDP F-series

หลังจากตัดแท็กแล้วเอกสารจะถูกตัด stopwords และทำการสเต็มมิ่ง ได้เป็นเซตของวลีที่สกัดได้จากเอกสารตัวอย่างนี้ ดังนี้

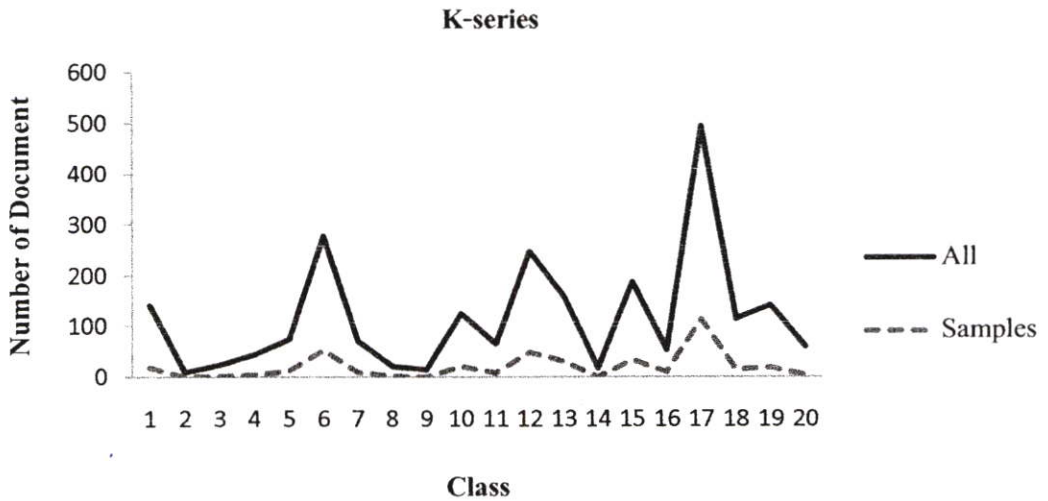
brought	edi implement	map problem	requir corpor chang
cal	edi readi	non trivi	requir surmount
capabl	exchang edi	paperfre system	small task
constitut	implem	perestroika	solution
corpor	integr edi	practic	theori
current busi practic	judgment	prepar	toughest obstacl technic
databas	known	problem	trad partner
difficult	liken	process	understood
direct result	maintain	program issu	widespread
docum	mak edi easier	proven	winmap address
edi	manageri	reluct	

ก.2 การสุ่มตัวอย่างเอกสารของชุดเอกสาร PDDP K-series

ชุดเอกสาร K-series มีเอกสารทั้งหมด 2,340 แบ่งเป็น 20 ประเภท ทำการสุ่มตัวอย่างแบบชั้นภูมิ โดยใช้อัตราส่วนในการเลือกเอกสารจากแต่ละชั้นภูมิเท่ากับ 0.05 รูปที่ ก.2 แสดงแผนภูมิแท่งของชั้นภูมิแต่ละประเภทเอกสาร โดยแกนนอนคือ ชั้นภูมิที่เป็นค่าต่างๆ และแกนตั้งคือจำนวนเอกสาร ที่มีค่านั้นๆ อยู่ โดยในที่นี้ได้เรียงลำดับจากชั้นภูมิที่มีจำนวนเอกสารน้อยสุดไปมากที่สุด หลังการสุ่มแล้ว ได้เอกสารตัวอย่างในแต่ละชั้นภูมิดังแสดงในตารางที่ ก.1 โดยประเภทที่มีเครื่องหมาย * คือประเภทเอกสารที่มีเอกสารจำนวนน้อย ซึ่งเราจะไม่นำมาใช้ในการทดลอง โดยจากเอกสารที่สุ่มได้ในตอนแรก 397 เอกสาร 20 ประเภทก็จะเหลือ 387 เอกสาร 13 ประเภท และสัดส่วนของจำนวนเอกสารแสดงในรูปที่ ก.3



รูปที่ ก.2 แสดงชั้นภูมิค่าและจำนวนเอกสารในแต่ละประเภทเอกสารของชุดเอกสาร K-series



รูปที่ ก.3 แสดงสัดส่วนระหว่างจำนวนเอกสารทั้งหมดกับเอกสารตัวอย่างในแต่ละประเภท ของชุดเอกสาร K-series

ตารางที่ ก.1 แสดงจำนวนเอกสารตัวอย่าง K-series แบ่งตามประเภทเอกสาร

ชื่อประเภท	จำนวนเอกสาร	
	ทั้งหมด	ตัวอย่าง
Business (B)	142	18
Entertainment (E)*	9	0
Art (Ea)*	24	1
Cable (Ec)*	44	4
Culture (Ecu)	74	12
Film (Ef)	278	52
Industry (Ei)	70	9
Media (Em)*	21	1
Multimedia (Emm)*	14	0
Music (Emu)	125	21
Online (Eo)	65	8
People (Ep)	248	48
Review (Er)	158	30
Stage (Es)*	18	0
Television (Et)	187	33
Variety (Ev)	54	10
Health (H)	494	113
Politics (P)	114	14
Sports (S)	141	19
Technology (T)*	60	4
รวม	2,340	397

ก.3 ชุดเอกสาร Reuters-top10

ชุดเอกสาร Reuters-top10 มีเอกสารที่อยู่ในรูปแบบ XML มีรายละเอียดและขั้นตอนการเตรียมเอกสารต่างๆ ดังนี้

ก.3.1 ตัวอย่างไฟล์เอกสาร Reuters-top10 ในรูปแบบ XML

รูปแบบของเอกสารก่อนการตัดแท็ก เป็นดังนี้

```
<?xml version="1.0" ?>
```

```
MESSAGE>USDA REPORTS EXPORT SALES ACTIVITY WASHINGTON, April 9 - The U.S.
```

```
Agriculture Department said private U.S. exporters reported sales of 200,000 tonnes of wheat to Jordan, 300,000 tonnes of soybean meal to Iraq and 100,000 tonnes of corn to Algeria. The wheat for Jordan includes 165,000 tonnes of hard red winter and 35,000 tonnes of soft red winter and is for delivery during the 1987/88 marketing year. The soybean meal sales to Iraq includes 180,000 tonnes for delivery during the 1986/87 season and 120,000 tonnes during the 1987/88 season, the department said. The 100,000 tonnes of corn sales to Algeria are for delivery during the 1986/87 season, it said. The marketing year for wheat begins June 1, corn September 1, and soybean meal October 1. Reuter</MESSAGE>
```

ก.3.2 ตัวอย่างไฟล์เอกสาร Reuters-top10 ที่ตัดแท็กออกแล้ว

```
usda reports export sales activity washington, april 9 - the u.s. agriculture department said private u.s. exporters reported sales of 200,000 tonnes of wheat to jordan, 300,000 tonnes of soybean meal to iraq and 100,000 tonnes of corn to algeria. the wheat for jordan includes 165,000 tonnes of hard red winter and 35,000 tonnes of soft red winter and is for delivery during the 1987/88 marketing year. the soybean meal sales to iraq includes 180,000 tonnes for delivery during the 1986/87 season and 120,000 tonnes during the 1987/88 season, the department said. the 100,000 tonnes of corn sales to algeria are for delivery during the 1986/87 season, it said. the marketing year for wheat begins june 1, corn september 1, and soybean meal october 1. reuter
```

ก.3.3 วิธีที่สกัดได้จากตัวอย่างไฟล์เอกสาร Reuters-top10

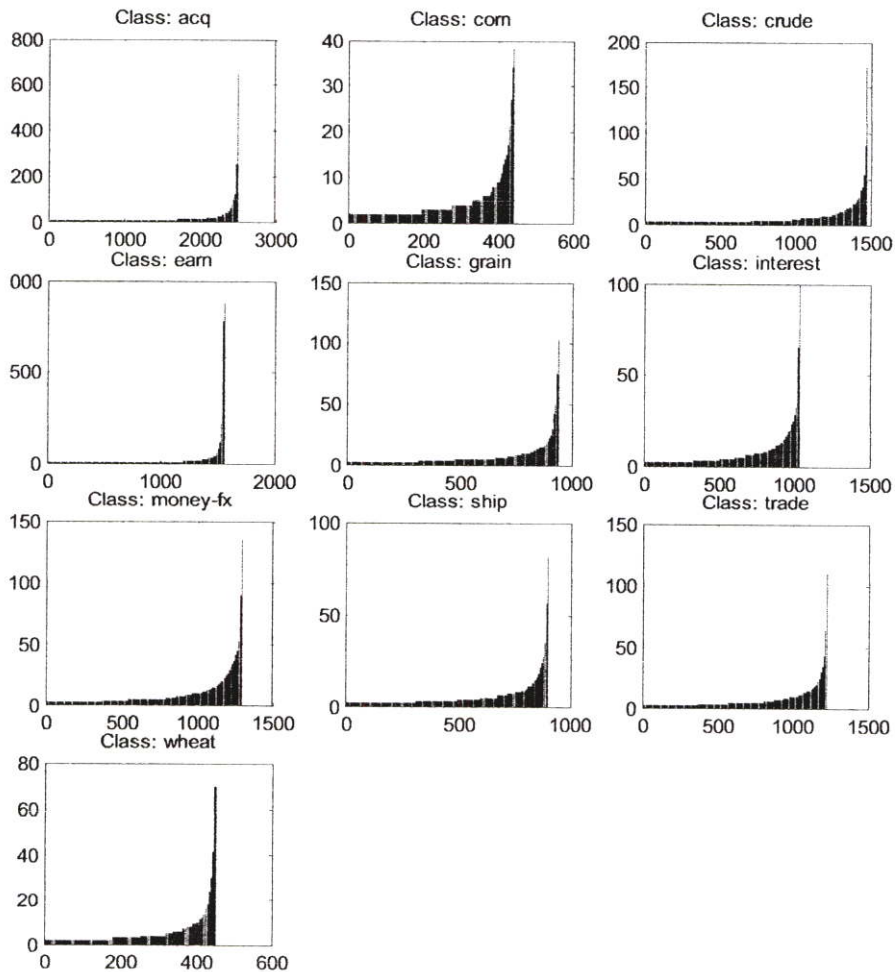
หลังจากตัดแท็กแล้วเอกสารจะถูกตัด stopwords และทำการสเต็มมิ่ง ได้เป็นเซตของวลีที่สกัดได้จากเอกสารตัวอย่างนี้ ดังนี้

agricultur depart said	export report sal	privat	soybean meal sal
algeria	hard red	reuter	tonn
april	iraq	said	usda report export sal activ
corn	iraq includ	season	washington

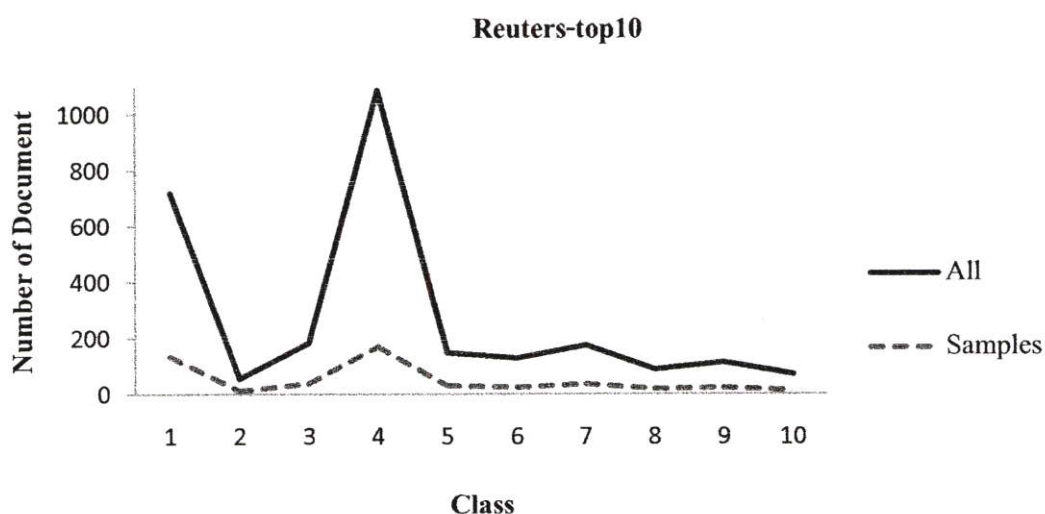
corn sal	jordan	soft red winter	wheat
corn septemb	jordan includ	soybean meal	wheat begin jun
deliveri	market year	soybean meal octob	winter
depart said			

ก.4 การสุ่มตัวอย่างเอกสารของชุดเอกสาร Reuters-top10

ชุดเอกสาร Reuters-top10 มีเอกสารทั้งหมด 2,770 แบ่งเป็น 10 ประเภท ทำการสุ่มตัวอย่างแบบชั้นภูมิ โดยใช้อัตราส่วนในการเลือกเอกสารจากแต่ละชั้นภูมิเท่ากับ 0.1 รูปที่ ก.4 แสดงแผนภูมิแท่งของชั้นภูมิแต่ละประเภทเอกสาร โดยแกนนอนคือ ชั้นภูมิที่เป็นค่าต่างๆ และแกนตั้งคือจำนวนเอกสาร ที่มีค่านั้นๆ อยู่ โดยในที่นี้ได้เรียงลำดับจากชั้นภูมิที่มีจำนวนเอกสารน้อยสุดไปมากที่สุด หลังการสุ่มแล้ว ได้เอกสารตัวอย่างในแต่ละชั้นภูมิดังแสดงในตารางที่ ก.2 ได้เอกสารตัวอย่างทั้งหมด 504 เอกสารจากทั้ง 10 ประเภท สัดส่วนของจำนวนเอกสารแสดงในรูปที่ ก.5



รูปที่ ก.4 แสดงชั้นภูมิค่าและจำนวนเอกสารในแต่ละประเภทเอกสารของชุดเอกสาร Reuters-top10



รูปที่ ก.5 แสดงสัดส่วนระหว่างจำนวนเอกสารทั้งหมดกับเอกสารตัวอย่างในแต่ละประเภทของชุดเอกสาร Reuters-top10

ตารางที่ ก.2 แสดงจำนวนเอกสารของชุดเอกสาร Reuters-top10 แบ่งตามประเภทเอกสาร

ชื่อประเภท	จำนวนเอกสาร	
	ทั้งหมด	ตัวอย่าง
acq	719	135
com	56	12
crude	184	38
earn	1,085	169
grain	148	31
interest	129	25
money-fx	176	37
ship	89	19
trade	113	24
wheat	71	14
รวม	2,770	504

๓.5 อัลกอริทึม Phrase-tree-based Feature Reduction

Algorithm: Phrase-tree-based Feature Reduction

Input: D – training set of documents; $d = d_1, d_2, \dots, d_n$

C – set of categories; $C = c_1, c_2, \dots, c_m$

L – set of label of each documents in D ; $L(d_i) \in C$

min_f – minimum frequency

Output: F : Set of Features

$F = \emptyset$

// Phrase I : Phrase-tree Construction & Local Reduction

For $i=1$ to m

$P = \emptyset$

For each document $d_j \in D$ that $L(d_j) = c_i$

PH = set of extracted phrases from d

$P = P \cup PH$

End

$P_tree_i = \text{construct_phrase_tree}(P)$ *//generate phrase-tree rom P*

$P_tree_i = \text{calculate_odds_ratio}(P_tree_i)$ *//calculate OR for each node in tree*

$P_tree_i = \text{reduce_phrase_tree}(P_tree_i, \text{min_f})$ *//Local Reduction*

End

// Phrase II : Phrase-tree Merging & Global Reduction

$M_tree = \emptyset$

For $i=1$ to m *//merge phrase-tree from each category*

$M_tree = \text{merge_tree}(M_tree, P_tree_i)$

End

$M_tree = \text{reduce_phrase_tree}(M_tree, \text{min_f})$ *//Global Reduction*

$F =$ all nodes from M_tree

ก.6 รายการคำหยุด (Stopwords) ที่ใช้ในการทดลอง

a	caption	get	kg	nobody	sm	using
ii	cc	gf	kh	none	sn	uy
about	cd	gg	ki	nonetheless	so	uz
above	cf	gh	km	noone	some	v
according	cg	gi	kn	nor	somehow	va
across	ch	gl	kp	not	someone	vc
actually	ci	gm	kr	nothing	something	ve
ad	ck	gmt	kw	now	sometime	very
adj	cl	gn	ky	nowhere	sometimes	vg
ae	click	go	kz	np	somewhere	vi
af	cm	gov	l	nr	sr	via
after	cn	gp	la	nu	st	vn
afterwards	co	gq	last	nz	still	vu
ag	co.	gr	later	o	stop	w
again	com	gs	latter	of	su	was
against	copy	gt	lb	off	such	wasn
ai	could	gu	lc	often	sv	wasn't
al	couldn	gw	least	om	sy	we
all	couldn't	gy	less	on	sz	we'd
almost	cr	h	let	once	t	we'll
alone	cs	had	let's	one	taking	we're
along	cu	has	li	one's	tc	we've
already	cv	hasn	like	only	td	web
also	cx	hasn't	likely	onto	ten	webpage
although	cy	have	lk	or	text	website
always	cz	haven	ll	org	tf	welcome
am	d	haven't	lr	other	tg	well
among	de	he	ls	others	test	were
amongst	did	he'd	lt	otherwise	th	weren
an	didn	he'll	ltd	our	than	weren't
and	didn't	he's	lu	ours	that	wf
another	dj	help	lv	ourselves	that'll	what
any	dk	hence	ly	out	that's	what'll

anyhow	dm	her	m	over	the	what's
anyone	do	here	ma	overall	their	whatever
anything	does	here's	made	own	them	when
anywhere	doesn	hereafter	make	p	themselves	whence
ao	doesn't	hereby	makes	pa	then	whenever
aq	don	herein	many	page	thence	where
ar	don't	hereupon	maybe	pe	there	whereafter
are	down	hers	mc	per	there'll	whereas
aren	during	herself	md	perhaps	there's	whereby
aren't	dz	him	me	pf	thereafter	wherein
around	e	himself	meantime	pg	thereby	whereupon
arpa	each	his	meanwhile	ph	therefore	wherever
as	ec	hk	mg	pk	therein	whether
at	edu	hm	mh	pl	thereupon	which
au	ee	hn	microsoft	pm	these	while
aw	eg	home	might	pn	they	whither
az	eh	homepage	mil	pr	they'd	who
b	eight	how	million	pt	they'll	who'd
ba	eighty	however	miss	pw	they're	who'll
bb	either	hr	mk	py	they've	who's
bd	else	ht	ml	q	thirty	whoever
be	elsewhere	htm	mm	qa	this	NULL
became	end	html	mn	r	those	whole
because	ending	http	mo	rather	though	whom
become	enough	hu	more	re	thousand	whomever
becomes	er	hundred	moreover	recent	three	whose
becoming	es	i	most	recently	through	why
been	et	i'd	mostly	reserved	throughout	will
before	etc	i'll	mp	ring	thru	with
beforehand	even	i'm	mq	ro	thus	within
begin	ever	i've	mr	ru	tj	without
beginning	every	i.e.	mrs	rw	tk	won
behind	everyone	id	ms	s	tm	won't
being	everything	ie	msie	sa	tn	would
below	everywhere	if	mt	same	to	wouldn

beside	except	il	mu	sb	together	wouldn't
besides	f	im	much	sc	too	ws
between	few	in	must	sd	toward	www
beyond	fi	inc	mv	se	towards	x
bf	fifty	inc.	mw	seem	tp	y
bg	find	indeed	mx	seemed	tr	ye
bh	first	information	my	seeming	trillion	yes
bi	five	instead	myself	seems	tt	yet
billion	fj	int	mz	seven	tv	you
bj	fk	into	n	seventy	tw	you'd
bm	fm	io	na	several	twenty	you'll
bn	fo	iq	namely	sg	two	you're
bo	for	ir	nc	sh	tz	you've
both	former	is	ne	she	u	your
br	formerly	isn	neither	she'd	ua	yours
bs	forty	isn't	net	she'll	ug	yourself
bt	found	it	netscape	she's	uk	yourselves
but	four	it's	never	should	um	yt
buy	fr	its	nevertheless	shouldn	under	yu
bv	free	itself	new	shouldn't	unless	z
bw	from	j	next	si	unlike	za
by	further	je	nf	since	unlikely	zm
bz	fx	jm	ng	site	until	zr
c	g	jo	ni	six	up	z
ca	ga	join	nine	sixty	upon	
can	gb	jp	ninety	sj	us	
can't	gd	k	nl	sk	use	
cannot	ge	kc	no	sl	used	

ภาคผนวก ข.

พีเจอรืของชุดเอกสารต่างๆ ที่ได้จากการทดลอง

ข.1 พีเจอรืที่ได้จากการเลือกโดยใช้ต้นไม้วลีของชุดเอกสาร F-series

พีเจอรืที่ได้จากการเลือกโดยใช้ต้นไม้วลีจากชุดการทดลองทั้ง 4 ชุดได้จำนวนพีเจอรืทั้งหมด 3,569 พีเจอรื ดังนี้

abbrevi	acknowledg	adopt standard ident	agenc mean
abfd	acom	advanc	agent
abil	acquir	advanc manufactur	aggreg
abl	acquisi	advanc manufactur system	aggress
abrasion	act	advanc manufactur technologi	agil
absenc	act mean	advantag	agil manufactur
absolut	action	advantag nasa	ago
abstract	action offic	advent	agre
ac	action plan	advers	agreem
ac polici	action program	advert	agricultur
academ	activ	advertis	agricultur servic
accept	activ involv	advic	agricultur worker
access	actual	advis	ahead
accid	ada	advisor	aid
accommod	adam	advoc	aid design
accomod	adapt	aerogel	aid manufactur
acompani	adapt layer	aerospac	aid manufactur engin
accomplish	add	affactweb	aim
accord	addition	affect	air
accordingli	address	affect affirm	alberta
account	adequ	affili	alert
account receiv	adjust	affirm	algorithm
accredit	administ	affirm action	alien
accredit standard	administr	affirm action plan	alleg
accredit standard committe	administr law	affirmative	alli
accru	administr law judg	affirmative ac	alloc
accur	administr procedur	affirmative ac polici	alloi
accuraci	admit	afford	allow
achiev	adop	afi	allow user
acid	adopt	afi cio	alpha
ack	adopt standard	agenc	alphabet

alter	applic	ask	audiovisu
altern	applic form	aspect	auditorium
aluminum	applic procedur	assault	august
amend	applic protocol	assembl	austin
america	applic system	assembli	authent
american	applications	assembli compon	author
american manufactur	applications ori	assembli compon requir	authorship
amount	appoint	assert	auto
amp	appreci	assess	autom
amsant	appris	asset	autom system
amus	approach	assign	automat
analog	appropri	assist	automot
analogi	approv	assoc	avail
analys	approxim	associ	avenu
analysi	april	assum	averag
analyt	arbitr	assur	avoid
analyt tool	arbitrag	astm	awai
analyz	arbor	asynchron	awar
anchorang	architectur	asynchron transfer	award
anecdote	archiv	asynchron transfer mod	axi
angel	area	atlanta	back
anim	area includ	atm	backbon
ann	argu	atm bas	background
ann arbor	argum	atm cell	bad
announc	aris	atm forum	bai
annual	arisen	atm network	bai area
annuiti	arizona	atm switch	balanc
ansi	arizona stat	atom	ball
answer	arizona stat univers	atp	ballot
antenna	arlington	attach	band
anti	arm	attack	bandwidth
anticip	armi	attempt	bank
antiqu	arrai	atten	banker
apart	arrang	attend	bar
apc	arriv	attorney	bar cod
appar	art	attorney fee	bargain
appeal	articl	attorney gener	bargain agreem
appeal rul	artist	attract	barrel
appear	artist work	attribut	barrel plat
appli	asc	audio	barrier
appli scienc	asid	audio tap	bas

bas client	bona fid	burden	car
bas network	bona fid occup	bureau	car analysi
basi	bond	burn	carbon
basic	book	burnish	card
basic fact	booklet	busi	career
battl	bor	busi administr	carefulli
baud	borrow	busi applic	camegi
baud rat	boston	busi cycle	camegi mellon
bear	bother	busi docum	camegi mellon univers
began	bottleneck	busi oper	carolina
begun	box	busi plan	carri
behalf	branch	busi process	carrier
behavior	brass	busi transac	cas
believ	break	butt	cas involv
bell	bridg	buyer	cas suggest
belong	brief	byte	cash
beneficiari	briefli	cabl	cash flow
benefit	bright	cabl televi	cassett
benefit payment	brilliant	cad	catalog
berkelei	bring	cad model	catalyst
bern	brittl	cad system	categori
bern conven	broad	cait	caught
best	broadband	cal	caus
best manufactur	broadcast	calendar	caution
best manufactur practic	broadli	california	cell
best manufactur practic survei	brochur	california civil	cellular
best practic	broken	california civil right	cellular telephon
better	broker	california civil right initi	censu
better repeat	brought	call	center
bid	browser	caller	centr
big	budget	caller number	central
bit	budget plan	cam	centuri
black	buffer	cam system	certain
blank	bui	canada	certain condition
block	build	candid	certain minimum
blow	build good	cap	certainli
board	build good credit	capabl	certif
boat	built	capac	certif requir
bob	bulk	capit	certifi
bodi	bulletin	capit fund	cfr
bona	bundl	captur	cfr part

chain	civil right act	combus	competi
challeng	civil right initi	comfort	competit
chanc	civil servic	command	competit advantag
chang	civilian	commenc	competit advantag nasa
channel	claim	comment	competitor
chapter	clarifi	commerc	compil
chapter appli	class	commerc integr	complain
charact	classic	commerc integr facil	complaint
characterist	classif	commerc technologi	complet
charg	classifi	commerci	complet applic
charl	claus	commerci offic	complet stat
chart	clean	commis	complex
cheaper	clear	commission	compli
check	clearanc	commit	complianc
chemic	clearinghous	committe	complianc program
chemistri	clearli	commod	complic
chen	clemson	common	compon
chicago	clemson univers	common databas	compon materi
chief	clerk	common languag	compon requir
child	client	common law	compos
children	client incorpor	commonli	composi
choic	clinton	commonli refer	composit
choos	clock	commun	compound
chosen	clos	commun colleg	comprehens
christoph	club	commun infrastructur	compress
chronolog	enc	commun servic	compris
chuck	coat	commun softwar	compulsori
chuck mclean	cod	communist	compulsori licens
cio	cold	compani	compuserv
cir	collabor	compani establish	comput
circl	collater	compani polici	computer
circuit	collec	compani progress	computer aid
circuit court	collect	company	computer aid manufactur
circular	collect bargain	company specif	computer aid manufactur
circumst	collect bargain agreem	compar	engin
cit	colleg	compar feder	conceiv
citi	collin	compar feder standard	concentr
citizen	color	compat	concept
civ	columbia	compens	conceptu
civil	column	compensatori	concern
civil right	combin	compet	conclud

conclusion	constitut	cool	court held
concurr	construc	cooper	court order
concurr engin	construct	cooper research	cover
condition	consult	coordin	cover public
conduct	consum	coordin council	coverag
conduct busi	consum credit	copi	cpu
conduct patent	consum group	copper	creat
conduct workplac	consum nation	copyright	creation
conduct workplac inspec	consum nation bank	copyright act	creativ
conduit	consump	copyright law	credit
confer	contact	copyright materi	credit card
confid	contact fcnb	copyright offic	credit manag
confidenti	contact milner	copyright owner	credit manag tip
configur	contact milner llnl	copyright propos	crew
confin	contain	copyright protec	criteria
conflict	contempl	copyright protect	critic
conform	content	copyright work	critical
confus	contentsbrbr	cor	crop
confusingli	context	cornell	cross
confusingli similar	conting	corner	cryogen
confusion	contin	corp	crystal
conges	contin improv	corpor	csl
congress	contract	correc	ct
congression	contract relat	correct	cultur
conjunc	contractor	correct action	currenc
connec	contractu	correctli	current
connect	contrari	correspond	curriculum
connect industri	contrast	corrosion	custom
connect industri repres	contribu	cosmet	custom produc
connecticut	contribut	cost	custom produc capabl
consensu	control	cost sav	custom satisfac
consent	control system	costli	custom servic
consequ	controll	council	cut
conserv	controversi	counsel	cut tool
consid	convei	counter	cwa
consider	conven	counti	cwa local
consist	conveni	countri	cwatx
consortia	conver	coupl	cyberspac
consortium	convert	cours	cycle
constant	convert debt	court	dai
constitu	convinc	court decision	daili

damag	deduct	detec	director form
danger	deem	detect	dirt
dat	deep	deterior	disa
data	defeat	determin	disabl
data bas	defect	develop	disabl act
data bit	defend	develop design	discard
data gener	defens	develop environ	discard polici
data interchang	defin	develop project	discharg
data link	defini	develop solution	disciplin
data manag	definit	devic	disclaim
data path	degrad	devic driver	disclos
data signal	degre	devot	disclosur
data structur	delai	di	disclosur docum
data transfer	delawar	diagnost	discount
data type	deliv	diamet	discov
databas	deliveri	diamond	discoveri
davi	demand	diamond tool	discretionari
david	demonstr	diamond turn	discrimin
day	deni	diamond turn machin	discrimin law
day haul	denial	dictat	discriminatori
dead	dental	dietitian	discus
deal	depart	differ	discuss
dealer	depend	differ type	disk
dealt	deploi	differenti	dismiss
dearborn	deploym	difficult	dispar
death	deposi	difficulti	displai
debat	deposit	digit	dispos
debt	depositori	digit commun	disput
debt instrum	depriv	digit equip	dissemin
dec	depth	digit equip corpor	distanc
decad	deriv	dimen	distinct
decemb	describ	dimension	distinctli
decid	descrip	diminish	distinguish
decision	descript	din	distribu
decisionmak	design	dip	distribut
declin	design methodologi	direc	district
decod	design methodologi group	direct	divers
decomposi	design process	directli	divid
decre	design research	directli involv	dividend
dedic	desir	director	division
deduc	detail	directori	doc

doctrin	earli	electron	encod
docum	earli stag	electron commerc	encount
dod	earlier	electron commerc integr	encourag
doe	earn	electron commerc integr facil	encourag stat
dol	earth	electron commerc technologi	encryption
dol mean	easi	electron data	energi
dollar	easier	electron data interchang	energi act
domain	easili	electron exchang	energi oak
domain logic	east	electron mail	energi oak ridg
domest	east avenu	electropl	energi system
domin	echo	elem	enforc
door	ecif	elig	engag
doubl	econom	elimin	engin
dout	economi	email	engin applic
download	edg	embodi	engin depart
downsid	edi	embrac	engin design
draft	edi implement	emerg	engin educ
draft intern	edi messag	emis	engin function
draft intern standard	edi program	emphasi	engin laboratori
drain	edi standard	emploi	engin societi
dramat	edi translat	employe	engin softwar
draw	edi user	employe mean	engin support
drawn	edi work	employe relat	engin tool
dream	edifact	employe train	engin tool kit
dri	edition	employer	engin toolkit
dril	educ	employer need	enhanc
driv	educ foundat	employer shall	enroll
driven	educ plan	employm	ensur
driver	eeo	employm agenc	enter
drop	eeoc	employm discrimin	enterpr
drop cell	effect	employm law	enterpr integr
drug	effect dat	employm offic	enterpr integr model
dry	effici	employm opportun	enterpr mod
drying	effort	employm opportun commis	enterpr represent
dual	einet	employm practic	entir
ductil	elect	employm relat	entiti
durabl	elect offici	employm servic	entitl
durat	electr	employm servic activ	entrepreneur
duti	electro	enabl	entri
dye	electro opt	enact	envelop
dynamic	electroless	encapsul	environ

environment	examin	expressli	favor
envision	exempl	ext	fax
enzym	exce	exten	fcnb
equal	excell	extend	fear
equal employm	excep	extend benefit	feasibl
equal employm opportun	excess	extens	featur
equal employm opportun	exchang	extent	feature
commis	exchang data	extern	februari
equal opportun	exclud	extra	feder
equilibrium	exclus	extrac	feder agenc
equip	exclus right	extrem	feder civilian
equip corpor	execu	fabric	feder contract
equip leas	execut	fac	feder copyright
equiti	execut order	facil	feder copyright protec
equiti particip	exempt	facilit	feder financi
equival	exerc	facsimil	feder financi assist
era	exercis	fact	feder fund
ergonom	exhaus	fact sheet	feder govern
erosion	exhaust	fact sheet osha	feder law
error	exhibit	factor	feder regist
escap	exist	factori	feder regul
especi	exist group	faculti	feder standard
essenc	expan	fail	federal
essenti	expand	failur	fee
establish	expect	fair	feel
estat	expedit	fairli	fell
estim	expens	fall	feng
eta	experi	fals	ferment
etch	experienc	famili	ferri
ethernet	experiment	famili member	fiber
ethnic	expert	familiar	fid
european	expir	famou	fid occup
evalu	explain	faq	field
event	explan	far	field offic
eventu	exploit	farm	fight
evid	explor	farm labor	figur
evolu	export	farm labor contractor	fil
evolv	expos	farmwork	fil amend
exact	exposur	fashion	fil fee
exactli	expres	fast	film
exagger	express	faster	filter

fin	forth	gari	graph
final	forum	gari woitena	graphic
financ	forward	gather	great
financi	foster	gaug	great deal
financi assist	foundat	gav	greater
finish	fourth	gear	greater invest
fir	fram	gender	greatest
firm	framework	gener	greatli
fiscal	franc	gener issu	greatli exagger
fiscal year	franchis	gener manag	green
fit	frank	gener principl	green paper
fix	fraud	gener provision	greet
fixat	fre	gener requir	grew
flat	frequenc	geograph	gross
flexibl	frequenc band	geometr	ground
flood	frequent	geometri	group
floor	frequent ask	georg	grow
florida	fri	georgia	growth
flow	friend	georgia tech	guarante
flow control	ftp	giant	guard
fluctuat	fuel	giv	guid
fmmla	full	given	guidelin
focu	full tim	glass	hair
focus	fulli	glob	half
follow	function	global	hand
food	fund	goal	handbook
foot	fund sourc	goe	handicap
footnot	fundament	gon	handl
forbid	furnac	good	happen
forc	furnish	good credit	happi
forecast	fusion	gopher	harass
foreign	futur	got	hard
foreign patent	futur atm	govern	hardwar
forese	futur atm network	govern agenc	harm
forget	futur contract	govern laborator	haul
form	futur work	government	hav
formal	gain	grad	hawaii
format	galaxi	graduat	hazard
formul	gam	grant	head
fort	gap	grant fund	header
fort collin	garden	grante	headquart

health	hopefulli	import	industri repres
health act	hospit	impos	industri research
health administr	host	imposs	industry
health car	hot	improp	inexpens
health plan	hotlin	improv	info
health program	hour	improv manufactur	inform
health standard	hous	improv product	infrastructur
hear	howard	inadequ	infrastructur compon
heard	hug	inch	infring
heat	human	incid	infring suit
heat treat	human interven	includ	ingredi
heavi	human resourc	includ addition	inherit
height	human resourc manag	includ school	initi
held	hundr	incom	injuri
hexagon	hypothet	incompat	ink
high	ibm	incomplet	innoc
high accuraci	iceimt	inconveni	innov
high degre	idaho	incorpor	input
high growth	idea	increas	input link
high perform	ideal	increasingli	input rat
high preci	ident	increasingli import	insert
high spe	identif	incur	insight
high speed	identifi	independ	inspec
high speed data	ieee	index	inspect
high speed network	igc	indian	inspector
high tech	ignor	indiana	inspir
higher	iii	indic	instal
higher level	ill	individu	install
highest	illinois	individu employ	instanc
highli	illustr	induc	institui
highli recommend	imag	industri	institut
highlight	imagin	industri associ	instruc
hir	immedi	industri classif	instrum
histor	immers	industri collabor	instrument
histori	impact	industri consortium	insur
hiv	impair	industri engin	intact
hol	implem	industri group	integr
hold	implement	industri need	integr circuit
holder	implement guidelin	industri particip	integr computer
holiday	implement plan	industri partner	integr computer aid
hop	implic	industri partnership	integr computer aid

manufactur	internet sit	item	jurisdic oper complet
integr division	internet support	iuliano	jurisdic oper complet stat
integr facil	internetwork	jack	jurren
integr model	interoper	jam	just
integr technologi	interpret	jan	justic
intellectu	interrupt	januari	justic depart
intellectu properti	interven	jlt	justif
intellectu properti right	interview	jlt uofrlaw	kansa
intellig	introduc	job	kansa technologi
inten	invalid	job applic	kbyte
inten discrimin	inven	job bank	keep
intend	invent	job descrip	kei
intens	inventor	job function	kentucki
intent	inventori	job open	kept
inter	invest	job opportun	kevin
interac	invest fund	job safeti	kevin jurren
interact	investig	job search	kevin waltz
interchang	investor	job servic	kind
interconnec	invit	job train	kit
interconnect	invit industri	job train partnership	know
interdisciplinari	invoice	job train partnership act	knowledg
interest	involuntari	joe	known
interfac	involv	joe milner	lab
interfac specif	involv copyright	john	label
interfer	ion	johnson	labor
interim	iowa	joint	labor contractor
intern	iowa stat	joint labor	labor manag
intern busi	iowa stat univers	joint labor manag	labor manag committe
intern commun	ipandc	joint labor manag committe	labor market
intern confer	ipande llnl	joint project	labor organ
intern copyright	ipo	jointli	labor regul
intern edi	ipr	jon	labor relat
intern revenu	iron	journal	labor union
intern revenu cod	island	judg	laboratori
intern standard	iso	judgment	laboratori invit
intern standard statu	iso standard	judici	laboratori invit industri
intern trad	isol	juli	laboratori multidisciplinari
intern trad administr	isp	juli parker	laboratori multidisciplinari
internation	issu	jun	capabl
internet	issu involv	jurisdic	lack
internet column	issu patent	jurisdic oper	lai

lak	left	livermor nation	machin part
lamin	legal	livermor nation laboratori	machin process
lan	legal entiti	livermor nation laboratori	machin tool
languag	legal protec	invit	machine
larg	legibl	livermor nation laboratori	madison
larg number	legisl	invit industri	magazin
large	legisl histori	llnl	magic
large scal	legitim	llnl directori	mail
larger	lend	llnl directori form	mail address
largest	lender	llnl disclaim	mail cod
laser	length	load	mail list
lat	letter	loan	mail order
latest	level	local	mailbox
lath	level packet	local employm	main
laundri	leverag	local govern	maintain
law	liabil	local law	mainten
law judg	liabl	local telephon	major
law relat	liaison	locat	mak
law requir	librari	lock	maker
law technologi	licens	log	manag
lawrenc	lif	logic	manag committe
lawrenc livermor	lif cycle	long	manag effort
lawrenc livermor nation	light	long term	manag mean
lawrenc livermor nation	lik	longer	manag system
laboratori	likelihood	look	manag tip
lawsuit	limit	loop	manner
lawyer	limit radio	loos	manual
layer	limit scop	los	manufactur
layout	lin	loss	manufactur applic
lbl	lin rat	lost	manufactur associ
lead	link	lot	manufactur autom
leader	liquid	loui	manufactur bas
leadership	list	louisiana	manufactur compani
learn	listen	low	manufactur data
leas	literari	low prioriti	manufactur educ
leaseback	literari work	lower	manufactur engin
leav	literatur	lower cycle	manufactur engin laboratori
lectur	litig	lower level	manufactur engin tool
lectur seri	littl	lowest	manufactur engin tool kit
led	liv	machin	manufactur enterpr
lee	livermor	machin connect	manufactur enterpr integr

manufactur equip	massiv	metal finish	modest
manufactur facil	match	metall	modif
manufactur network	match fund	meter	modifi
manufactur oper	materi	method	modifi march
manufactur part	materi engin	methodologi	modul
manufactur practic	materi scienc	methodologi group	moistur
manufactur practic survei	matrix	metk	molecular
manufactur problem	matter	metrologi	molecular model
manufactur process	matur	metropolitan	moment
manufactur program	maximum	metropolitan area	monci
manufactur research	mbp	mexico	monica
manufactur resourc	mclean	mhz	monitor
manufactur resourc data	mean	michael	monopoli
manufactur sandia	mean occup	michigan	month
manufactur sandia nation	mean servic	mid	month period
manufactur scienc	meant	middl	monthli
manufactur simul	measur	migrant	monthli payment
manufactur support	mechan	mik	mortgag
manufactur system	mechan engin	mik iuliano	motion
manufactur system integr	media	militari	motion pictur
manufactur technologi	medic	milner	mount
manufacturing	medicin	milner lnl	mov
map	medium	mind	movem
map problem	meet	minim	msfw
mar	mel	minimum	msid
march	mel thrust	minimum cost	mtu
margin	mellon	minneapolis	multi
marietta	mellon univers	minor	multi media
marietta energi	member	minut	multi servic
marietta energi system	membership	mississippi	multidisciplinari
mark	memori	mistak	multidisciplinari capabl
market	mental	misunderstand	multimedia
market innov	mention	mit	multimedia applic
market maker	mer	mix	multimedia commun
marketplac	mer idea	mobil	multipl
martin	merger	mobil commun	multiplex
martin marietta	merit	mod	mural
martin marietta energi	messag	model	mural project
martin marietta energi system	messag specif	modem	murrai
mask	met	modern	music
mask work	metal	modern produc	mutual

mutual fund	network bas	numer control	oper complet
nacha	network interfac	nurs	oper complet stat
nam	network provid	nyper	oper cost
nasa	network resourc	nyper public	opinion
nation	network resourc manag	oak	oppon
nation bank	network testb	oak ridg	opportun
nation institut	neural	oak ridg center	opportun commis
nation laboratori	neutral	oak ridg nation	oppos
nation laboratori invit	newli	oak ridg nation laboratori	opt
nation laboratori invit industri	newslett	obes	optic
nation law	newslett summar	object	optim
nation origin	newslett summar court	object ori	optimist
nation research	newslett summar court	object orient	optimum
nation research council	decision	oblig	option
nation scienc	nickel	observ	order
nation scienc foundat	nist	obstacel	order describ
nation technolog	nist laboratori	obtain	orderli
nation technolog univers	nist manufactur	obviou	ordinari
nationwid	nist system	obvious	ordinarili
nativ	nist system laboratori	occasion	oregon
nativ american	nistir	occup	organ
natl	nixon	occup safeti	organ interest
natur	nois	occur	organiz
navi	non	octob	ori
navig	non profit	odd	orient
ncsl	non profit organ	ofcep	origin
nearest	nondiscrimin	offer	origin work
nearli	nonexclus	offic	ornl
necessari	nonprofit	offic receiv	osha
necessarili	normal	offic requir	osha approv
necessit	north	offici	osha conduct
need	north american	offset	osha conduct workplac
neg	northwest	ohio	osha conduct workplac inspec
neg consequ	notic	oil	osha inspec
negoti	notif	old	osha program
neighbor	notifi	older	osha standard
neil	novel	omit	otc
neil christoph	novemb	onli	outer
nell	nuclear	onlin	outgo
netcom	number	open	outlin
network	numer	oper	outperform

output	particip	permiss	plan oper
output link	particular	permit	plan oper cost
output port	particular fact	perpetu	plan system
outreach	particularli	person	plant
outsid	partner	person asset	plastic
overhead	partner agreem	person characterist	plat
overrid	partnership	person contact	platform
overview	partnership act	personnel	plaza
owner	pass	personnel manag	plu
ownership	past	perspect	point
ownership position	patent	pertain	point diamond
oxford	patent applic	peys	polari
oxford univers	patent attornei	peys act	polari plat
oxid	patent avail	phas	polic
pac	patent law	philadelphia	polic offic
packag	patent offic	phon	polic
packet	patent protec	phon call	polish
packet siz	patent right	phon compani	polit
packet switch	patent search	phonograph	polymer
packet switch network	path	photocopi	polytechn
pag	pattern	photograph	polytechn institut
pai	paul	phras	poor
paid	payload	physical	popul
paint	payment	physical harm	popular
pair	pc	physical locat	port
panel	peddl	physical scienc	portabl
paper	pen	pick	portfolio
paperfre	penalti	pictur	portion
paperfre system	pencil	piec	posit
paragraph	pend	pilot	position
parallax	peopl	pilot implement	posses
parallax educ	perceiv	pin	possess
parallel	percent	pip	possibl
paramet	percentag	pip system	possibl applic
paramount	percep	pirat	possibl infring
parent	perform	plac	possibli
park	perform comput	placem	post
parker	period	plain	post offic
part	period mean	plan	postscript
parti	perman	plan applic	potenti
partial	permis	plan goal	pound

poverti	prim	product	protect
pow	primari	product data	protect trademark
power	primarili	product data manag	protect worker
power suppli	princip	product lin	protocol
practic	principi	product qualiti	prototyp
practic survei	print	profession	prov
pre	print version	professor	proven
precaution	prior	profil	provid
preced	prioriti	profit	provid certain
preci	privat	profit margin	provid equal
precis	privat employer	profit organ	provision
precision	privat firm	profit shar	provision govern
precision engin	privat investor	program	pub
precision machin	privat placem	program budget	public
precision manufactur	privat sector	program budget plan	public accommod
predict	privileg	program develop	public employe
predominantli	probabl	program highlight	public employm
prefer	problem	program provid	public law
preferenti	problem aris	programm	public librari
preferenti treatment	proce	progress	public offer
preliminari	procedur	prohibi	public perform
premis	proceed	prohibit	public school
premium	process	projec	public servic
prepar	process control	project	publicli
prescrib	process develop	project includ	publish
present	process equip	promin	puerto
preserv	process involv	promot	puerto rico
presid	process plan	promotion	pump
presid clinton	process plan applic	prompt	punish
presid johnson	process plan system	proof	pur
presid nixon	process specif	propag	purchas
presidenti	process specif languag	proper	purchas option
press	process worker	properli	purchas order
press releas	processor	properti	purpos
pressur	procur	properti law	pursu
presump	produc	properti right	pursuant
pretti	produc capabl	propos	put
prevent	produc facil	proprietary	qualif
previou	produc manag	proprietorship	qualifi
previous	produc manag system	prospect	qualifi person
pric	produc part	protec	qualiti

qualiti control	reader	regist trademark	repair
qualiti manag	readi	registr	repeal
qualiti manag effort	readili	regul	repeat
quantiti	readili access	regul concern	replac
quarterli	real	regul specif	replenish
question	real estat	regular	report
question aris	real tim	regularli	repres
queue	real world	regulatori	represent
quick	realiti	rehabilit	reprint
quickli	realiz	rehabilit act	reproduc
quit	realli	reinstat	republican
quot	realm	reject	reput
quota	ream	rel	request
rac	reason	relat	request edi
racial	reason accommod	relat activ	request form
racism	reason attorney	relat board	requir
rack	reason attorney fee	relationship	requir elem
rack plat	reason opportun	releas	requir employer
radio	recall	relev	research
radio spectrum	receipt	reli	research center
rag	receiv	reliabl	research council
rai	recipi	relianc	research group
rais	recogn	relief	research laborator
rais capit	recogni	reliev	research program
rang	recommend	religi	resembl
rapid	record	religion	reserv
rapid access	recordkeep	reluct	resid
rapid prototyp	recov	remain	resist
rapid respons	recoveri	remedi	resolu
rapid respons manufactur	recreat	rememb	resolv
rapid tool	recruit	remitt	resourc
rapidli	reduc	remot	resourc data
rar	refer	remot manufactur	resourc manag
rat	referr	remot procedur	respect
rat fell	refund	remov	respond
ratio	refus	remov soil	respond neg
reach	regard	render	respons
react	regardless	renew	respons manufactur
reaction	region	renssela	rest
read	region offic	rent	restor
read copyright	regist	rental	restric

restrict	robot	sba field	segment
result	robust	sba field offic	segreg
retail	rock	sbic	sel
retain	rol	scal	selec
retali	rom	scar	select
retir	room	scenario	self
retrain	root	sched	sell
return	roughli	schedul	semant
rev	round	schem	semiconductor
reveal	rout	scholarship	seminar
revenu	router	school	send
revenu cod	routin	scienc	sender
review	royalti	scienc foundat	senior
revis	rrm	scientif	sens
revision	rt	scientist	sensit
revok	rul	scop	sensor
revolu	run	scor	sent
revolution	rural	screen	separ
rhod	saf	search	sept
rich	safeti	searchabl	septemb
richmond	sai	season	seq
richmond journal	said	season farmwork	sequenc
rico	sal	sec	seri
rid	salari	second	serial
ridg	sampl	secondari	serial commun
ridg center	sampl issu	secondari offer	serv
ridg nation	san	secret	server
ridg nation laboratori	sandi	secretari	servic
right	sandia	secretari mean	servic activ
right act	sandia nation	secretary	servic agenc
right confer	sandia nation laboratori	secretary treasur	servic mark
right initi	santa	secretary treasur cwa	servic mean
rigid	santa monica	secretary treasur cwa local	servic mean servic
rins	satellit	section	servic perform
rippl	satisfac	sector	servic system
ris	satisfactori	secur	set
risk	satisfi	secur act	set forth
riski	sav	see	set job
road	sav account	seed	set job safeti
road map	saw	seek	settl
robert	sba	seen	settle

setup	silicon	small firm	specialist
sever	silicon vallei	smaller	specif
sex	sima	smaller institu	specif languag
sexual	sima program	sme	specif technologi
sexual discrimin	similar	sme educ	specifi
shall	similar mark	sme educ foundat	spectrum
shall appli	similar right	smok	specul
shap	similarli	smooth	speech
shar	simpl	social	speed
shar memori	simpl credit	social secur	speed data
sharehold	simpl credit manag	social secur act	speed network
shaw	simpl credit manag tip	societi	spend
shaw feng	simpler	soft	spin
sheet	simplest	softwar	sponsor
sheet osha	simpli	softwar applic	spread
shelf	simpli phon	softwar tool	spur
shift	simplifi	softwar vendor	stabil
ship	simul	soil	staf
shipment	simultan	sol	staff
shop	singl	sold	stag
shop floor	singl physical	solicit	stai
short	single	solid	stainless
short form	single point	solubl	stak
short term	single point diamond	solution	stand
short term debt	sit	solv	standard
shorter	situ	somewhat	standard committe
shot	situat	song	standard cover
show	siz	soon	standard develop
shown	skill	sophist	standard ident
sic	skin	sort	standard industri
sick	slic	sought	standard iso
sid	slightli	sound	standard number
sign	slim	sound record	standard requir
signal	slogan	sourc	standard requir employer
signal process	slop	south	standard set
signal rt	slow	south carolina	standard statu
signatur	slower	southern	standard titl
signific	small	spac	stanford
significantli	small busi	spe	stanford univers
silica	small busi administr	speak	stanlei
silica aerogel	small compani	special	start

stat	street	suit	system
stat agenc	street address	suitabl	system develop
stat agenc mean	strength	sum	system integr
stat cod	strengthen	summar	system integr division
stat employm	stress	summar court	system laboratori
stat employm servic	strik	summar court decision	system research
stat govern	strip	summari	system research center
stat job	striv	sun	systemat
stat law	strong	sup	tabl
stat patent	strong technologi	superintend	tactic
stat plan	strong technologi bas	supervi	tailor
stat program	strongli	supervisor	tak
stat unemploy	struck	supp	taken
stat unemploy compens	structur	supplementari	takeov
stat univers	student	suppli	talk
stat university	studi	suppli chain	tangibl
statem	sub	supplier	tank
station	sub project	support	tap
statist	subchapt	support disciplin	target
statu	subdivid	support servic	task
statut	subject	suppos	tax
statutori	subject matter	supra	tax advantag
statutori protec	submis	suprem	tax benefit
steal	submit	suprem court	tax credit
steam	subpart	sur	tb
steel	subscrib	surcharg	tcp
stem	subscrip	surfac	teacher
step	subsequ	surfac area	team
stereolithographi	substanc	surpris	tech
steven	substanti	surround	techn
stewart	substitut	survei	technic
stick	subtl	surviv	technic approach
stock	subtract	suscept	technic assist
stock offer	success	suspen	technic resourc
stor	success stem	suspend	techniqu
storag	successfulli	sustain	technolog
stori	suddenli	sweden	technolog need
straightforward	sue	switch	technolog univers
strateg	suffici	switch fabric	technolog
strategi	sugges	switch network	technolog bas
streamlin	suggest	symbol	technolog develop

technologî need	ti	trademark	trying
technologî program	tight	trademark law	turbochannel
technologî transfer	tim	trademark offic	turn
technologî transfer offic	time	trademark right	turn machin
tel	timet	trademark suit	turn machin tool
telecommun	tip	tradition	tutori
telecommun industri	tir	traffic	twist
telephon	titanium	train	type
telephon directori	titl	train administr	typewrit
telephon lin	titl iii	train cours	typical
telephon number	toc	train institut	ucrl
telephoni	today	train partnership	ucrl tb
televi	toler	train partnership act	ultim
tell	toll	train program	unabl
temperatur	toll fre	transac	unauthor
temperatur measur	tomorrow	transfer	unawar
temporari	took	transfer mod	unchang
temporarili	tool	transfer offic	uncontrol
tend	tool appropri	translat	underli
tent	tool avail	translat softwar	understand
term	tool data	transmis	understood
term debt	tool kit	transmit	undertak
termin	toolkit	transmitt	underwrit
territori	topic	transport	undesir
testb	topologi	trap	unemploy
tetherless	total	travel	unemploy compens
texa	total domain	treasur	unequ
thank	total qualiti	treasur cwa	unfair
theater	total qualiti manag	treasur cwa local	unfair competi
theori	total qualiti manag effort	treasuri	unfortun
thermal	touch	treat	unifi
thick	track	trati	unifi process
thing	track record	treatment	unifi process specif
think	trad	tree	unifi process specif languag
third	trad administr	tremend	union
thorough	trad associ	trend	uniqu
thoroughli	trad commis	tri	unit
thought	trad develop	troubl	unit stat
threat	trad partner	true	unit stat cod
throughput	trad partner agreem	trust	univers
thrust	trad secret	try	univers librari

university	vast	wagner peys	wireless network
unix	vehicl	wagner peys act	wisconsin
unlaw	vendor	wai	wish
unlaw employm	ventur	wait	wit
unlaw employm practic	ventur capit	waiver	withdrew
unrel	verifi	walk	withhold
unusu	version	waltz	woitena
uofrlaw	versu	want	woman
updat	veteran	war	women
updat mar	viabl	warrant	wong
updat sched	vic	wash	word
upgrad	victim	washington	work
upheld	video	wast	work capit
upper	vietnam	water	work dai
urg	vietnam era	waterloo	work forc
urich	view	wav	work plan
url	vigor	wax	work schedul
usa	vii	weak	worker
usa ornl	viii	weaken	worker employ
usabl	violat	weapon	worker mean
usag	violenc	wear	workforc
usdoc	virgin	webmast	workplac
user	virgin island	week	workplac inspec
user data	virginia	weekli	workshop
user network	virtual	weight	worksit
usual	virtual circuit	welfar	workstat
utah	virtual enterpr	west	world
util	virtual enterpr mod	western	world leader
vacat	virtual librari	whit	world wid
vacuum	virtual manufactur	whit hous	worldwid
valid	virtual realiti	whit paper	worri
vallei	vision	wholli	worth
valu	visit	wid	wrap
valuabl	visual	wide	writ
value	vocat	wider	written
vapor	voic	widespread	written descrip
vari	volum	wil	year
variabl	voluntari	window	year ago
variat	vot	wir	year mean
varieti	wag	wireless	year plan
variou	wagner	wireless commun	year plan goal

yield	zinc
york	zip
young	zip cod
youth	
zemo	
zero	

ข.2 พีเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วิของชุดเอกสาร J-series

พีเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วิจากชุดการทดลองทั้ง 4 ชุดได้จำนวนพีเจอร์ 3,488 พีเจอร์ ดังนี้

abandon	achiev	administr comput	affili
abil	acknowledg	admis	affirm
abl	acoust	admit	affirm action
abroad	acquir	adop	affirm action polici
absenc	acquisi	adopt	affirm action program
absent	act	advanc	afford
absolut	action	advanc composit	african
absorp	action polici	advanc comput	african american
abstrac	action program	advanc manufactur	ag
abstract	activ	advanc materi	ag discrimin
academ	activist	advanc system	agenc
academia	actual	advanc technologi	agent
acceler	actuat	advantag	aggreg
accept	ad	advers	agil
access	adapt	advert	agre
accid	add	advertis	agreem
accommod	addit	advic	agricultur
accompani	addition	advis	aid
accompani measur	address	advisor	aid design
accomplish	address includ	advisori	aim
accord	adequ	adviz	air
accordingli	adher	advoc	air forc
account	adhes	aeronaut	airborn
account receiv	adhesion	aerospac	aircraft
accur	adjust	affair	airport
accuraci	administ	affid	alamo
accus	administr	affect	alamo nation

alamo nation laborator	apparel manufactur	artist	automobil
alert	appeal	artist work	automot
algebra	appear	asian	autonom
algorithm	appell	asian american	avail
align	appelle	ask	avenu
alleg	appli	ask question	averag
allegedli	applic	aspect	aviat
allen	applic develop	assembl	avoid
alli	applic form	assembli	awai
allianc	applic program	asser	awalt
alloc	applic softwar	assert	awar
alloi	applic system	assess	award
allow	appoint	asset	award program
alston	approach	assign	axi
alter	appropri	assist	back
altern	approv	associ	background
aluminum	approxim	assum	backup
amend	april	assump	bad
america	aqueou	assur	balanc
american	arbitr	atlanta	band
amount	arc	atm	bank
amplifi	architect	atmosph	bank busi
analys	architectur	atom	bank loan
analysi	architectur design	attach	banker
analysi techniqu	archiv	attain	bar
analyst	area	attempt	barbara
analyt	area includ	atten	bargain
analyz	area network	attend	bargain agreem
anc	argu	attomei	bargain unit
andersen	argum	attomei fee	barrier
andersen consult	aris	attract	bas
anim	arkansa	attribut	bas approach
anisotrop	aros	audienc	bas data
announc	arrai	audio	bas design
annual	arrang	audit	bas industri
answer	art	augment	bas manufactur
anticip	articl	august	bas materi
api	artifici	authent	bas system
appar	artifici intellig	author	bas technologi
apparatu	artifici neural	autom	basi
apparel	artifici neural network	automat	basic

batch	bipartisan	busi capit	capit fund
batteri	bipartisan commis	busi capit corpor	capit invest
beam	black	busi cas	capit market
bear	block	busi develop	captur
began	board	busi electron	car
begun	bodi	busi enterpr	carbon
behalf	bond	busi environ	card
behavior	book	busi financ	career
behaviour	boost	busi innov	carolina
belief	borrow	busi manag	carri
believ	boston	busi need	carrier
bell	boundari	busi object	cas
benchmark	boundari layer	busi owner	cas law
benefici	box	busi plan	cas studi
benefit	brak	busi process	cash
benign	braz	busi process engin	cash flow
berkeley	breach	busi system	cast
berlin	break	busi transac	cast process
bern	bridg	button	cat
bern conven	brief	buyer	catalog
bern union	bring	cabl	catalysi
best	broad	cad	catalyst
best practic	broad rang	cal	catalyt
better	broadcast	calcul	categori
bial	broaden	calendar	caus
bid	broader	calibr	ceil
big	broadli	calif	ceil commis
bilater	brought	california	cell
bind	brown	call	cellular
bio	browser	cam	cent
biochem	bu	cambridg	center
biolog	budget	campu	centr
biolog process	bui	canada	central
biolog system	build	canadian	centrifug
biologi	built	candid	centuri
biomass	bulk	cap	ceram
biomed	burden	capabl	cert
biomed applic	burn	capac	certain
biomimet	bush	capit	certainli
bioremedi	busi	capit corpor	certif
biotechnologi	busi administr	capit format	certifi

chain	circumst	columbia	complementari
chair	cisba	com	complet
chairman	cit	combin	completion
challeng	citi	combus	complex
chamber	citizen	combus engin	complex system
chanc	civil	comfort	compli
chang	civil engin	comm	complianc
channel	civil right act	command	complic
chapter	civil servic	commenc	compon
charact	civilian	comment	compos
character	claim	commerc	composi
characterist	class	commerc depart	composit
charg	classic	commerc requir	composit materi
charl	classif	commerci	composit structur
chart	classroom	commerci exploit	compound
check	claus	commis	comprehens
chemic	clean	commit	compres
chemic engin	clear	committe	compress
chemic process	clearli	common	compris
chemic reaction	client	common law	comprom
chemic vapor	client server	commonli	compulsori
chemic vapor deposi	climat	commun	compulsori licens
chemistri	clinton	commun network	comput
chen	clos	commun system	comput environ
chicago	cluster	commun technologi	comput system
chief	coat	compact	concentr
children	cod	compani	concept
china	coeffici	compar	conceptu
chip	coher	comparison	concern
choic	coincid	compat	concern patent
choos	cold	compel	conclud
chos	collabor	compens	conclusion
chosen	collater	compet	concurr
cim	collec	competi	concurr design
cim model	collect	competit	concurr engin
cim system	collect bargain	competit advantag	concurr engin environ
cir	collect bargain agreem	competitor	condition
circuit	colleg	compil	conduc
circuit board	collision	complain	conduct
circuit judg	colloid	complaint	confer
circul	color	complem	confid

confidenti	contamin	corpor	crystal growth
configur	contempl	correc	crystallin
conflict	contend	correct	ctiv
conform	content	correl	ctiv week
conform test	context	correspond	cultur
confront	conting	corrosion	current
cong	continuu	cost	curriculum
congress	continuu improv	cost effect	custom
congression	contract	cost sav	custom order
congression intent	contract claim	cost shar	custom servic
conjunc	contractor	costli	cut
connec	contractu	council	cycle
connect	contrari	counsel	cycle engin
consciou	contribu	count	cycle tim
consent	contribut	counter	dai
consequ	control	counti	daili
conserv	control technologi	countri	damag
consid	controll	coupl	dan
consider	conven	cours	daniel
consist	conver	court	dat
consolid	convert	court conclud	dat problem
consortia	cool	court decision	data
consortium	cooper	court grant	data acquisi
constant	cooper research	court held	data bas
constitu	cooper treati	coven	data collec
constitut	coordin	cover	data exchang
constrain	copi	coverag	data integr
constraint	copper	crada	data interchang
constru	copyright	creat	data manag
construc	copyright act	creat job	data model
construc materi	copyright conven	creation	data secur
construc process	copyright law	creativ	data structur
construct	copyright offic	credit	data transfer
consult	copyright owner	credit card	databas
consult group	copyright phillip	creditor	databas technologi
consult servic	copyright phillip publish	crimin	davi
consum	copyright protec	crisi	davi publish
consum product	copyright relat	criteria	davi publish compani
consump	copyright work	critic	david
contact	cor	cross	deadlin
contain	corp	crystal	deal

dean	depart	develop organ	disciplin
debat	department	develop propriitari	disciplinari
debentur	depend	develop propriitari research	disciplinari action
debit	deploi	develop requir	disclaim
debt	deploym	develop tim	disclos
debug	deposi	devic	disclosur
dec	deposit	devis	discount
decad	depreci	diagnost	discov
decemb	depriv	diagram	discoveri
decentralis	dept	dialogu	discret
decid	depth	dictat	discretion
decis	deputi	die	discrimin
decision	deregul	differ	discriminatori
decision analysi	deriv	differ level	discus
decision mak	describ	differ type	discuss
decision support	descrip	differenti	diseas
decision support system	design	difficult	disk
decreas	design decision	diffract	dismiss
dedic	design environ	diffusion	dispatch
deduc	design method	diffusion bond	displai
deem	design patent	digit	dispos
deep	design process	digit equip	disposi
defeat	design tool	digit equip corp	disput
defect	desir	digit signal	dissemin
defend	desk	digit signal process	dissent
defens	desktop	dimen	distinct
defer	despit	dimension	distinguish
defin	destroi	diminish	distribu
defini	destruct	diod	distribut
deform	destruct evalu	dioxid	distribut databas
degrad	detail	direc	distribut environ
degre	detec	direct	distribut manufactur
delai	detect	direct invest	district
deliv	determin	directli	district court
deliveri	develop	director	district court grant
demand	develop advanc	directori	divers
democrat	develop agreem	disabl	diversif
demonstr	develop cycle	disabl benefit	diversifi
deni	develop design	disadvantag	division
denomin	develop method	disappear	do
densiti	develop methodologi	discharg	doctor

doctrin	easili	electron system	energi
docum	east	electron technologi	energi effici
document	econom	elem	energi storag
doe	econom develop	elig	enforc
dol	econom growth	elimin	eng
dollar	econom impact	email	engag
domain	economi	email address	engin
domest	ectodai	embed	engin applic
donat	ed	embodi	engin center
door	edg	emerg	engin design
doubt	edi	emerg applic	engin environ
downstream	edit	emis	engin materi
downward	edition	emori	engin program
downward pressur	editor	emphas	engin system
dozen	editor continu	emphasi	engin technologi
draft	educ	emphasi ad	english
drag	educ servic	empir	engr
dramat	edward	emploi	enhanc
draw	eeo	employe	enorm
drawn	effect	employe benefit	ensur
driv	effici	employe receiv	ensur equal
driven	effort	employe right	entail
driver	elast	employe right continu	enter
drop	elec	employe work	enterpr
drug	elector	employer	enterpr integr
dsp	electr	employ	enterpr model
dual	electr power	employ act	enterpr wid
due	electro	employ contract	entir
duk	electro optic	employ practic	entiti
durabl	electromagnet	employ relationship	entitl
durat	electromechan	empow	entrepreneur
duti	electron	enabl	entrepreneuri
dynamic	electron beam	enabl technologi	entri
earli	electron catalog	enact	environ
earlier	electron circuit	encapsul	environment
earliest	electron data	encompass	environment benign
earn	electron data interchang	encourag	environment friendli
earth	electron devic	encryption	environment impact
eas	electron packag	end	environment protec
easi	electron payment	endheadingbr	environment remedi
easier	electron product	endors	environment technologi

envisag	exam	explain	fall
enzym	examin	explan	famili
epoxi	exampl	explicitli	familiar
equal	exampl includ	exploit	famou
equal employm	exce	explor	far
equal opportun	exceed	explos	fashion
equat	excell	export	fast
equilibrium	excep	expos	faster
equip	exception	exposur	fat
equip corp	excerpt	expres	fatigu
equip design	excess	express	fault
equit	exchang	express written	fault toler
equiti	exchang transac	express written permis	favor
equiti financ	exchang valu	expressli	fax
equiti invest	excim	exten	feasibl
equival	exclud	exten partnership	featur
era	exclus	exten servic	februari
ergonom	exclus right	extend	fed
error	execu	extens	feder
escap	execu system	extent	feder agenc
especi	execut	extern	feder court
esprit	execut order	extra	feder fund
essenti	execut summari	extrac	feder govern
establish	exerc	extract	feder labor
estat	exercis	extrem	feder law
estim	exist	extrusion	feder program
etch	exist applic	fabric	fee
ethic	exist busi	fabric process	feedback
ethnic	exist system	fac	feel
europ	exit	facet	fellowship
european	expan	facil	femal
european industri	expand	facilit	fewer
european market	expand busi	fact	fiber
evalu	expect	factor	fiber optic
evalu techniqu	expenditur	factori	fibr
event	expens	faculti	fiduciari
eventu	experi	fail	field
evid	experienc	failur	fifth
evolu	experiment	failur analysi	figur
evolutionari	expert	fair	figur show
evolv	expir	faith	fil

fil dat	flux	fred	genet engin
fil patent	fly	freedom	geograph
film	focal	french	geometri
filter	focu	frequenc	georg
fin	focus	frequent	georgia
fin particl	follow	friction	germani
final	food	friendli	get
final product	food process	fring	giv
financ	footnot	fring benefit	given
financ program	forc	frustrat	glass
financi	forecast	fuel	glass ceil
financi institu	foreign	fulfill	glass ceil commis
financi intermediari	foreign applic	fuller	global
financi market	foreign bank	fulli	global electron
financi servic	foreign countri	fulli integr	global electron commerc
financi support	foreign patent	function	global market
find	form	function requir	go
finish	formal	fund	goal
finish product	formal member	fund sourc	goe
finit	format	fund transfer	good
finit elem	formul	fundament	good faith
fir	forrest	furnish	gorrell
firm	forrest research	furthermor	govern
fiscal	forth	fusion	govern agenc
fiscal year	fortun	futur	government
fit	forum	fuzzi	governor
fix	forward	ga	grac
fix asset	foster	gain	grac period
flat	found	gam	grad
flat panel	foundat	gap	gradient
flexibl	fourth	gas	gradual
flexibl manufactur	fraction	gather	graduat
flight	fractur	gatt	graduat student
float	fractur mechan	gav	grant
floor	fram	gear	grant summari
floor control	framework	gen	grant summari judgment
florida	franc	gender	graphic
flow	franchis	gener	great
fluid	francisco	gener object	greater
fluid dynamic	frank	gener process	greatest
fluid flow	fraud	genet	greatli

greg	help	host	import
grievanc	helpdesk	hostil	import rol
grievanc procedur	hierarch	hot	impos
grind	high	hour	improv
gross	high densiti	hous	improv flexibl
ground	high frequenc	hub	improv product
group	high level	human	improv understand
groupwar	high perform	human activ	impur
grow	high perform comput	human factor	inadvert
growth	high power	human resourc	incent
guarante	high prioriti	human resourc manag	incid
guid	high puriti	hundr	includ
guidanc	high qualiti	hybrid	includ analysi
guidelin	high resolu	hydrocarbon	includ applic
half	high risk	ibm	includ cad
hamilton	high school	ic	includ chemic
hand	high technologi	idea	includ design
handicap	high temperatur	ideal	includ electron
handl	high volum	ident	includ metal
handl system	higher	identif	incom
happen	higher educ	identifi	inconsist
har	highest	ieee	incorpor
hard	highli	ignor	increas
hardwar	highli effici	iii	increasingli
harmon	highli integr	iitf	increasingli import
hav	highlight	ill	increment
hazard	highwai	illinoi	incur
hazard wast	hill	illinoi stat	indebted
head	hir	illinoi stat univers	independ
headquart	hispan	illustr	index
health	histor	ilstu	indic
health car	histori	imag	individu
hear	hoc	immedi	individu compani
heart	hold	impact	individu employe
heat	holder	impair	individu right
heat flux	holograph	implant	induc
heat transfer	hong	implem	industri
heat treatment	hong kong	implement	industri develop
heavi	hop	implement plan	industri engin
heavi metal	horizont	impli	industri firm
held	hospit	implic	industri includ

industri need	insul	intern conven	java
industri partner	insur	intern journal	jet
industri partnership	integr	internation	jit
industri partnership offic	integr circuit	internet	job
industri product	integr design	internet commerc	job creat
industri properti	integr manufactur system	internet commun	job creation
industri standard	integr model	internet market	job list
inequ	integr multi	interop	job market
inexpens	integr process	interpret	john
infer	integr servic	interv	johnson
influen	integr system	interven	join
info	integr technologi	interview	joint
inform	integr tool	introduc	journal
informix	intel	inven	judg
infosi	intellectu	invent	judgment
infrar	intellectu properti	inventor	judici
infrastructur	intellectu properti law	inventori	juli
infring	intellectu properti right	inventori control	jun
ing	intellectualproperti	inventori manag	jurisdic
inher	intellig	inventorship	just
initi	intellig materi	invers	justic
initi capit	inten	invest	justifi
injec	intend	invest strategi	kanban
injuri	intens	investig	kansa
innov	intensifi	investor	keep
innov technologi	intent	invit	kei
inorgan	inter	invoic	ken
input	inter ctiv	involv	kept
inquiri	inter ctiv week	ion	kernel
insid	interac	ipr	kim
insight	interact	irrelev	kind
inspec	interchang	island	kinemat
instal	interconnec	iso	kinet
install	interest	isol	kingdom
instanc	interest parti	issu	know
institui	interfac	item	knowledg
institut	interfer	jam	knowledg bas
instruc	intermedi	jan	known
instrum	intermediari	januari	kong
instrument	intern	japan	lab
insuffici	intern bank	japanes	labor

labor law	lean	list	low volum
labor manag	learn	literari	lower
labor organ	leas	literari work	lower cost
labor practic	leav	literatur	lower rat
labor relat	led	lithographi	lynx
labor relat act	lee	litig	maarten
labor relat board	left	littl	machin
laboratori	legaci	liv	machin system
lack	legaci system	livermor	machineri
laid	legal	livermor nation	machinist
lak	legal advic	livermor nation laboratori	macintosh
lamin	legal protec	llp	magazin
lan	legisl	lo	magnet
land	legitim	lo alamo	magnitud
languag	lend	lo alamo nation	mail
lanl	lender	lo alamo nation laboratori	mail comment
larg	letter	load	main
larg scal	level	loan	mainfram
larger	level model	local	mainli
largest	leverag	local govern	maintain
laser	lewi	locat	mainten
laser bas	lewi research	logic	major
laser process	liabil	logist	major bank
lat	liaison	logo	major intern
latest	librari	london	mak
launch	licenc	long	maker
law	licens	long term	mal
law firm	license	long term capit	man
law provid	lif	longer	manag
lawrenc	lif cycle	longer term	manag need
lawsuit	lift	look	manag practic
lawyer	light	loop	manag strategi
layer	light weight	loss	manag system
layout	lightweight	lost	manipul
layout design	likelihood	lot	manner
lead	limit	lot siz	manual
lead edg	lin	loui	manufactur
lead tim	linear	low	manufactur applic
leader	link	low cost	manufactur compani
leadership	linkag	low level	manufactur cost
leakag	liquid	low temperatur	manufactur data

manufactur design	mathemat	mesh	mislead
manufactur engin	matric	messag	mission
manufactur enterpr	matrix	met	mitig
manufactur exten	matrix composit	metabol	mix
manufactur exten partnership	matter	metal	mixtur
manufactur improv	matur	metal cast	mobil
manufactur industri	maxim	metal form	mocvd
manufactur oper	maximum	metal matrix	mod
manufactur requir	mean	metal matrix composit	modal
manufactur softwar	meant	metall	model
manufactur system engin	measur	metallurg	model allow
manufactur system integr	measur techniqu	metallurgi	model bas
map	mechan	meter	model data
march	mechan engin	method	model techniqu
margin	mechan properti	methodologi	modem
mari	mechan system	metr	moder
marin	mechatron	metric	modern
maritim	media	metrologi	modif
mark	medic	mfg	modifi
market	medium	mi	modul
market plac	medium siz	michael	modular
market relat	medium siz manufactur	michigan	mold
marketplac	medium term	mid	molecul
mass	meet	middl	molecular
mass produc	melt	middlewar	monei
mass transfer	mem	migrat	monitor
massachusett	member	mik	monitor system
master	member countri	militari	monolith
mastercard	membership	mill	monopoli
mastercard intern	membran	min	month
match	memori	mind	mortgag
match fund	men	miner	mosaic
materi	ment	miner process	motion
materi develop	ment contract	miniatur	motiv
materi flow	mention	minim	motor
materi handl	menu	minimis	mountain
materi process	mep	minimum	mov
materi properti	mer	minist	movem
materi research	merchant	minor	mrp
materi scienc	merg	minor group	multi
materi technologi	merit	misconduct	multi nation

multichip	neural	nrtw	opportun
multidisciplinari	neural network	nrwldf	oppos
multifunc	neutral	nsf	opposit
multilater	neutron	nuclear	optic
multilay	new	number	optic commun
multimedia	newli	numer	optic fiber
multin	newslett	numer control	optim
multipl	newspap	numer method	optim problem
music	newsstand	numer model	optim techniqu
nam	ng	oak	optimum
narrow	nich	object	option
nasa	nickel	object orient	optoelectron
nasa lewi	nist	oblig	or
nat	nlra	observ	oracl
nation	nlr	obstacl	oracl corp
nation engin	nod	obtain	oral
nation institut	nois	obviou	order
nation labor	non	occasion	order process
nation labor relat	non destruct	occup	organ
nation labor relat act	non destruct evalu	occur	organ repres
nation labor relat board	non equilibrium	ocean	organ solvent
nation laboratori	non profit	octob	organis
nation origin	nondestruct	offer	organiz
nation scienc	nondestruct evalu	offic	orient
nation scienc foundat	nondestruct test	offici	orient model
natur	nonlinear	ohio	origin
natur languag	nonlinear dynamic	oil	origin work
navi	nonprofit	old	outlin
navig	normal	omit	output
near	north	on	output start
nearli	north america	ongo	outreach
necessari	north american	onlin	outreach program
necessarili	not	open	outsid
need	notic	open system	outstand
neglect	notif	oper	overcom
neglig	notifi	oper condition	overlook
negoti	notwithstanding	oper cost	oversea
net	novel	oper manag	overview
netherland	novelti	oper system	ow
network	novemb	opinion	own
network servic	novo	opm	owner

ownership	patent applic	pharmaceut	polit
oxid	patent attornei	phas	poll
ozon	patent cooper	phenomena	pollut
pac	patent cooper treati	phenomenon	pollution
pacif	patent inven	phillip	polymer
packag	patent issu	phillip logo	polymer materi
pag	patent law	phillip publish	pool
pag copyright	patent offic	phillip team	poor
pai	patent process	phon	poorli
paid	patent protec	photograph	pop
paint	patent relat	photon	popul
panel	path	phras	port
paper	pattern	physic	portabl
paper present	pattern maker	physical	portfolio
paradigm	pattern recogni	physical properti	portion
paragraph	paul	physician	pos
parallel	payment	physiolog	posit
parallel comput	pc	physiologi	position
paramet	pdt	pilot	possibl
pari	penalti	pilot project	post
pari conven	pend	pirat	potenti
park	pennsylvania	plac	poverti
part	peopl	placem	powder
parti	percent	plai	power
partial	percentag	plain	power electron
particip	percentag point	plaintiff	power gener
particl	perform	plaintiff claim	practic
particul	perform comput	plan	practition
particular	perform measur	planar	pre
particular emphasi	period	plant	preced
particularli	perman	plant oper	precis
partner	permis	plasma	precision
partnership	permiss	plastic	preclud
partnership offic	permit	platform	predic
pass	person	plu	predict
passag	personnel	plug	prefer
passeng	personnel manag	point	preferenti
passiv	perspect	pointer	preferenti treatment
past	persuas	polar	perform
past year	pertain	polic	prejudic
patent	petroleum	policymak	preliminari

premis	process	profession servic	proven
premium	process bas	professor	provid
prepar	process data	profil	provid electron
prerequisit	process design	profit	provid financ
prescrib	process develop	program	provid insight
presenc	process engin	program languag	provid integr
present	process includ	program manag	provid loan
preserv	process industri	program provid	provid user
presid	process manag	programm	provision
presid clinton	process manufactur	progress	pto
press	process model	prohibi	pub
pressur	process monitor	prohibit	public
presum	process plan	project	public kei
prevail	process simul	project cost	public polici
preven	process system	project manag	public sector
prevent	process techniqu	project select	publicli
previou	process technologi	promis	publish
previou section	process wast	promot	publish compani
previous	processor	promot electron	pulp
pric	procur	promotion	pur
prim	produc	prompt	purchas
prim minist	produc activ	proof	purif
primari	produc control	propag	puriti
primarili	produc environ	proper	purport
princip	produc manag	properli	purpos
principl	produc plan	properti	pursu
print	produc process	properti law	pursuant
prior	produc schedul	properti protec	push
prioriti	produc system	properti right	put
privaci	product	propos	qualif
privat	product announc	proposi	qualifi
privat network	product cost	proprietary	qualit
privat sector	product data	proprietary research	qualiti
privileg	product data model	propul	qualiti assur
pro	product design	prosecu	qualiti control
proactiv	product lif	prospect	qualiti improv
probabl	product lif cycle	protec	qualiti manag
problem	product model	protect	quantifi
proce	product qualiti	protocol	quantit
procedur	profes	prototyp	quantiti
proceed	profession	prov	que

queri	real tim	regularli	reproduc
question	realist	regulatori	reproduc right
quick	realiti	regulatori framework	republican
quick respons	realiz	reinforc	request
quickli	realli	reinstat	requir
quit	reason	reinvest	research
quot	receiv	reiter	research activ
quota	reces	reject	research area
rac	recipi	rel	research center
racial	recogn	relat	research group
racism	recogni	relat act	research institut
radar	recommend	relat board	research issu
radiat	reconfigur	relat data	reserv
radio	record	relat technologi	resid
rai	recours	relationship	residu
rai smock	recov	releas	resign
railroad	recoveri	relev	resin
rais	recruit	reli	resist
rais capit	recycl	reliabl	resolu
ralli	red	relianc	resolv
random	redesign	relief	resort
rang	redistribu	remain	resourc
rank	reduc	remand	resourc center
rapid	reduc cost	remark	resourc librari
rapid prototyp	reduc energi	remedi	resourc manag
rapidli	reengin	remot	respect
rapidli chang	refer	remov	respond
rat	referr	remov process	respons
ratio	refin	render	rest
ration	reflect	renew	restor
raw	reform	rental	restric
raw materi	refrain	reorgan	restrict
reach	refus	rep	result
reaction	regard	repair	retail
reactiv	regardless	replac	retain
reactor	regim	report	retain job
read	region	report examin	retali
readi	regist	repositori	retir
readili	registr	repres	retriev
real	regul	represent	retrofit
real estat	regular	reprint	return

reus	safer	section provid	servic
reusabl	safeti	sector	servic compani
revenu	sai	secur	servic group
revers	said	seed	servic mark
revers discrimin	sal	seek	servic organ
review	salari	seen	servic provid
revis	sampl	segment	servic sector
richard	san	sel	session
ridg	san francisco	selec	set
right	sap	select	set forth
right act	satisfac	self	settl
right confer	satisfi	sell	setup
right continu	sav	seller	sever
right issu	sba	semant	sex
rigor	sba loan	sematech	shall
ris	sbir	semi	shap
risk	scal	semiconductor	shar
risk associ	scan	semiconductor devic	shear
rival	scatter	semiconductor manufactur	sheet
road	schedul	seminar	shelf
robert	schem	senat	shell
robot	school	send	shift
robust	scienc	senior	ship
rol	scienc foundat	senior circuit	shop
roland	scientif	senior circuit judg	shop floor
rom	scientist	senior manag	shop floor control
room	scop	sens	short
ross	scopi	sensit	short term
rossi	scopic	sensit analysi	shorter
rossi alston	scrap	sensor	show
rotat	screen	sensor arrai	shown
round	script	sensor technologi	shut
rout	sea	sent	sign
routin	seamless	separ	signal
royalti	search	septemb	signal process
rue	sec	seq	signatur
rul	second	sequenc	signific
rul bas	secondari	seri	significantli
run	secret	serv	silicon
rural	secretari	server	similar
saf	section	server architectur	similarli

simpl	softwar modul	speed	stereotyp
simpli	softwar product	spend	stev
simplif	softwar system	spin	stev bell
simplifi	softwar tool	spirit	steven
simul	sol	sponsor	stimul
simul method	sold	spooft	stochast
simul model	solder	spot	stock
simul tool	solid	sprai	stock quot
simultan	solid phas	spread	stop
singl	solid stat	spreadsheet	stor
sinter	solidif	sru	storag
sit	solidif process	stabil	stori
situ	solution	stabl	strain
situat	solv	staf	strateg
sixth	solvent	staff	strategi
siz	somewhat	stag	stream
siz manufactur	sonic	stai	streamlin
skill	soon	stand	street
slight	sophist	standard	strength
slow	sort	start	strengthen
small	sought	stat	stress
small busi administr	sound	stat court	strict
small busi capit	sourc	stat district	strik
small busi develop	south	stat district court	striv
smaller	spac	stat govern	strong
smart	span	stat law	strongli
smart materi	spatial	stat patent	structur
smart structur	speak	stat personnel	structur design
sme	speci	stat personnel manag	structur materi
smgr	special	stat univers	struggl
smith	specialist	statem	student
smock	specialti	station	studi
social	specif	statist	su
societ	specif industri	statu	sub
societi	specif situat	statut	subject
soft	specifi	statutori	submis
softwar	spectral	statutori right	submit
softwar applic	spectroscopi	statutori schem	subordin
softwar develop	spectrum	steel	subsequ
softwar engin	specul	steelwork	subsidi
softwar implement	speech	step	subsidiari

substanc	surfac modif	team	textil worker
substanti	survei	teamster	thank
substitut	surveill	tech	them
substrat	surviv	technic	theoret
subsystem	sustain	technic assist	theori
success	sweet	technic support	therapeut
successfulli	switch	techniqu	thermal
sue	switzerland	technolog	thermal analysi
suffer	sybas	technolog	thermal conduct
suffici	symposium	technolog applic	thermodynam
suggest	synthesi	technolog area	thermoplast
suit	synthesiz	technolog bas	thermoset
suitabl	synthetic	technolog corpor	thing
sum	system	technolog develop	think
summar	system analysi	technolog includ	thoma
summari	system analyst	technolog program	thought
summari judg	system approach	technolog provid	thousand
summari judgment	system bas	technolog requir	threat
sun	system develop	tel	threaten
sun system	system engin	telecom	throughput
supercomput	system group	telecommun	ti
superconduct	system includ	teleoper	tightli
superconductor	system integr	telephon	tim
superior	system solution	telephon number	tim limit
superplast	systemat	televi	tim scal
superplast form	tabl	tell	tion
supersed	tactic	temperatur	tip
supervi	tailor	temporari	tissu
supervisor	tak	tend	titl
supp	taken	tension	titl vii
suppli	talk	ter	tiv
suppli chain	tangibl	term	tm
supplier	tap	term capit	today
support	target	termin	told
support develop	tariff	termin employe	toler
support servic	task	terrestri	took
support system	tax	test	tool
support tool	taxat	testabl	topic
suprem	tcp	testb	toronto
sur	teach	texa	total
surfac	teacher	textil	total number

total qualiti	truli	unit kingdom	vapor
toxic	trust	unit oper	vapor deposi
tqm	try	unit stat	vari
trac	tub	unit stat court	variabl
traceabl	turbin	unit stat district	variat
track	turbul	unit stat district court	varieti
track record	turn	unit stat patent	variou
trad	turnov	unit steelwork	vector
trad mark	tutori	univers	vehicl
trad relat	twic	univers commun	veloc
trad secret	twin	univers copyright	vendor
trademark	type	univers copyright conven	ventur
trademark applic	typical	unix	ventur capit
trademark offic	ultim	unlaw	verif
tradition	unabl	unpaid	verifi
tradition manufactur	unauthor	unwant	version
traffic	uncertainti	updat	vertic
train	undergradu	upgrad	vessel
transac	underground	upstream	veteran
transac onlin	underli	urban	vibrat
transac process	underpin	urg	vic
transfer	understand	urgent	vic presid
transfer model	understood	url	video
transform	undertak	us	view
transition	undertaken	usabl	vii
translat	underwai	usag	violat
transmis	unfair	user	violet
transport	unfair labor	user interfac	virtual
transport phenomena	unfair labor practic	usual	virtual enterpr
travel	unfortun	util	virtual realiti
treasuri	unifi	vacat	visa
treat	uniform	vacuum	vision
treati	union	valid	vision system
treatment	union disciplin	vallei	visit
treatment process	union effort	valu	visual
trend	union membership	valu ad	vital
trial	union rul	valu chain	voic
trip	union secur	valuabl	voltag
trp	union shop	valuat	volum
truck	uniqu	valv	volum produc
true	unit	van	voluntari

voluntarily	weight	withdraw	worldwid
vortic	weld	woman	worth
vot	went	women	worth not
wafer	west	wood	writ
wag	western	word	written
wai	whichev	work	written permis
wait	whit	work capit	wrong
wall	whit mal	work condition	wto
wang	whit men	work environ	www
want	wholesal	work group	year
war	wid	work program	year patent
wareh	wid rang	work relationship	yield
warehous	wid varieti	worker	york
wari	widespread	worker compens	young
washington	wil	worker local	zdnnet
wast	william	workflow	zdnnet logo
wast dispos	willing	workforc	zdnnet
wast stream	win	workplac	zero
water	wind	workshop	ziff
wav	window	workstat	ziff davi
wav process	wipo	world	ziff davi publish
wavelength	wir	world class	ziff davi publish compani
weapon	wireless	world intellectu	
wear	wireless commun	world intellectu properti	
weather	wish	world market	
week	wit	world wid	

ข.3 ฟังเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลีของชุดเอกสาร K-series

ฟังเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลีจากชุดการทดลองทั้ง 4 ชุดได้จำนวนฟังเจอร์ 3,731 ฟังเจอร์ ดังนี้

abandon	abnorm	academi award	accid victim
abc	abor	accent	accident
abc ti	abroad	accept	accompani
abdomen	absenc	access	account
abil	abus	access award	accumul
abl	academi	accid	accus

achiev	affect	aliv	analyst
achiev award	affili	alleg	analyst said
acknowledg	africa	allegedli	anderson
acquir	african	allen	andrew
acquisi	ag	allerg	andrew hind
act	ag associ	allerg reaction	andrew morton
act rock	ag associ diseas	allergi	aneurysm
action	ag bas	allergi try	aneurysm risk
activ	ag bas rat	allergi try wash	angel
actor	ag group	alli	angel counti
actress	ag relat	allow	angel lo
acut	ag relat deficit	alomar	angel lo angel
ad	ag year	alongsid	angel trial
adam	agenc	alter	angela
adapt	agent	altern	angioplasti
add	agent help	alto	angioplasti arteri
addic	agent help crohn	altogeth	angioplasti arteri clear
addition	agent help crohn diseas	alzheim	anim
address	aggress	alzheim brain	ann
adelaid	ago	alzheim brain damag	ann miller
adequ	ago averag	alzheim diseas	anni
adher	agre	amaz	anniversari
adjust	agreem	amend	announc
administ	agricultur	america	annual
administr	ahead	america onlin	annual meet
admir	aid	american	answer
admis	aim	american journal	anti
admis criteria	air	american leagu	antibiot
admit	air forc	american leagu championship	antibiot avoid
adult	airbag	american leagu championship	antibiot avoid ulcer
adult demograph	airbag injuri	seri	antibodi
adulthood	airwai	american medic	anticip
advanc	alan	american medic associ	antidepress
advers	albert	american societi	antigen
advers affect	album	ami	anxieti
advertis	album guid	amicu	anybodi
advis	album guid seri	amicu award	anymor
advisori	alcohol	ampl	apart
advisori committe	alert	amus	apologi
advoc	alert broaden	analysi	appar
affair	alex	analysi appear	appeal

appear	assn	avoid ulcer	batter
appl	associ	awai	batteri
appl spokeswoman	associ diseas	await	bean
applic	associ produc	awar	bear
approach	associ professor	award	beat
appropri	atkinson	award gala	beatl
approv	atkinson comedi	award go	beauti
approv irradi	atlant	award present	bed
approxim	atlanta	award win	beef
approxim half	attach	babi	began
archerd	attack	babi born	begin
archiv	attempt	bacall	begun
area	atten	back	behalf
argentin	attend	bacteria	behav
argu	attitud	bad	behav badli
aris	attornei	badli	behavior
arkansa	attract	ballad	behavior risk
arlington	auction	balloon	belief
arm	audienc	baltimor	believ
armi	aug	baltimor oriol	belt
armi archerd	august	band	belushi
arrang	australia	bank	benefit
arrest	australian	bar	benign
arriv	author	barbara	benz
art	author conclud	barbra	bernard
art direc	author not	barrier	best
art scienc	author sai	bart	best friend
arteri	author stat	bas	best friend wed
arteri clear	auto	bas rat	best tim
articl	autoimmun	baseball	bestsell
artist	automak	baseball gam	bethesda
ashlei	automak seen	baseball playoff	better
asia	automak seen post	basi	beverli
ask	automak seen post lukewarm	basic	beverli hill
assault	automak seen post lukewarm	bath	bid
assembl	profit	bath seat	big
assess	avail	bath seat risk	bigger
assign	averag	bath seat risk infant	biggest
assist	averag ag	bath seat risk infant liv	billi
assist director	avoid	bathroom	bind
assist professor	avoid drug	bathub	biochemistri

biolog	bowman	brown	campbell
biologi	bowman grai	browser	campbell endear
biologist	bowman grai school	buckl	campbell endear perform
biopsi	box	budget	canada
birth	box offic	buena	canadian
birthdai	brad	buena vista	canal
bishop	brain	bueno	cancer
bit	brain cell	bueno air	cancer cell
biz wir	brain damag	bui	cancer patient
black	brain diseas	build	cancer risk
blam	brain surgeri	built	candid
block	brain surgeri help	burden	cann
blockag	brain surgeri help parkinson	burn	cap
blond	brain tissu	busi	cap town
blood	branch	caa	capabl
blood pressur	brand	cabl	capac
blow	break	cabl network	capit
blue	breakthrough	cableac	capitol
blur	breast	cableac award	capitol hill
blur deliv	breast cancer	cag	captur
blur deliv mix	breast cancer risk	cal	car
blur deliv mix set	breath	calcul	car burden
bmi	brian	calendar	car cost
board	bridg	calendar item	car crash
bob	brief	calif	car worker
bodi	bring	california	card
bodi mass	brit	call	career
bodyguard	britain	calor	caregiv
boi	british	calor intak	carei
bon	british medic	calor intak prolong	carlo
bond	british medic journal	calor intak prolong lif	carol
book	bro	calor restric	carolina
boost	broad	calor restrict	carri
born	broadcast	calor restrict diet	cas
boss	broaden	calori	cast
boston	bronchoscop	cam	cat
bought	bronchoscop spread	camera	catch
bovin	bronchoscop spread tuberculosi	camerawork	categori
bovin spongiform	brook	camp	catherin
bow	brother	campaign	cathet
bowel	brought	campaign fund	cathol

caus	charact	chronic	coincid
caus breast	character	chronic ill	cold
caution	characterist	chrysler	coli
cb	charg	chrysler corp	colicki
cd	charl	church	colicki babi
cdc	charli	cigarett	colin
cdc nation	charlton	cigarett smok	collabor
cdc offici	charlton heston	cinema	colleagu
cecer	charm	circl	colleagu interview
cedar	chart	circul	colleagu point
cedar sinai	chas	cit	colleagu stat
celebr	chas ami	citi	collec
cell	cheer	civil	collect
cell immun	chemic	clai	collect prais
cell receptor	chemotherapi	claim	colleg
cell slow	chest	class	color
cell slow parkinson	chicago	classic	colorado
cellular	chicago hop	clean	columbia
cancel	chief	clear	columbia univers
cent	chief execut	clearli	column
center	chief film	cleveland	columnist
centr	chief film critic	cleveland indian	com
central	chil	clinic	combat
centuri	child	clinic trial	combin
ceo	childhood	clinton	comedi
ceremoni	children	clinton meat	comfort
certain	china	clinton meat recall	comic
certainli	chines	clinton meat recall run	command
cessat	chip	clos	commerc
chad	choic	closet	commerci
chad ogea	cholesterol	cloud	commerci evad
chain	choos	clout	commis
chairman	choreographi	club	commit
chairman gat	chos	clue	committe
challeng	chri	coach	common
champion	chri petrikin	coast	commun
championship	christian	coceain	compact
championship seri	christian soldier	cod	compact disc
chanc	christin	coffe	compani
chang	christma	coffe pric	compani said
channel	christoph	cognit	companion

companion dodi	contamin	cours	cyber
compar	contemporari	court	cyber summari
comparison	contend	cover	cynthia
competi	content	coverag	cynthia littleton
competit	contest	cow	dad
competitor	continu	cow link	dai
complain	contracep	cox	daili
complet	contracept	cpr	daili liv
complex	contract	cpr lesson	daili varieti
compli	contrari	craft	daili varieti chief
complianc	contrast	crash	daili varieti chief film
complic	contribu	craven	daili varieti chief film critic
compos	contribut	creat	daili varieti highlight
compound	control	creativ	daili varieti highlight octob
comput	control group	credibl	daili varieti senior
concentr	controversi	credit	daili varieti senior columnist
concep	conven	credit includ	damag
concept	convic	creutzfeldt	damag caus
concern	convict	creutzfeldt jakob	dan
concert	convinc	creutzfeldt jakob diseas	dan cox
conclud	cop	crimin	danazol
condemn	cop land	criteria	danc
condition	copi	critic	danger
conduct	cor	crohn	daniel
confer	cord	crohn diseas	danni
confidenti	comer	cross	dario
confirm	coronari	crowd	dark
congress	coronari arteri	cruis	dat
congression	coronari heart	cult	data
connecticut	coronari heart diseas	cultur	datelin
consecut	corp	cultur dish	daughter
consensu	corpor	cultur showbiz	davi
consid	cosbi	cultur showbiz event	david
consider	cosmet	cup	daytim
consist	cost	curb	dead
consolid	costum	curb teen	deaf
constant	costum design	curb teen smok	deal
consult	count	current	dean
consum	counti	current issu	dean goodman
contact	countri	current studi	death
contain	coupl	cut	death occur

death rat	demonstr	dick	dividend
debat	deni	die	division
debut	denni	diet	division seri
debut album	depart	dietari	divorc
dec	depend	differ	dna
decad	depres	difficult	dna repair
decemb	depress	difficulti	do
decid	derek	digit	doctor
decision	derek ellei	dinner	docum
declin	deriv	direc	documentari
decoi	dermatologist	direct	dodi
decontamin	describ	directli	dolbi
decor	deserv	director	dollar
decreas	design	director gener	domest
dedic	despit	director oliv	domin
dee	destroi	director oliv ston	donna
deep	destruc	disabl	door
defeat	detail	disappoint	dopa
defend	detect	disc	dos
defend world	deterior	discourag	dosag
defend world seri	determin	discov	doubl
defiantli	detroit	discoveri	doubt
defici	detroit piston	discrep	doug
deficit	develop	discus	dougla
defin	devic	discuss	downward
definit	devito	diseas	dr
degen	di	diseas associ	draft
degre	diabet	diseas caus	drama
delai	diabet risk	diseas control	dramat
delai school	diabet risk identifi	diseas occur	draw
delight	diagnos	diseas relat	dream
deliv	diagnosi	diseas transmis	dress
deliv mix	diagnost	dish	drew
deliv mix set	dialogu	disnei	drew carei
dell	diamond	disord	drift
dell corp	dian	distant	drink
demand	dian ladd	distribu	driv
demeanor	diana	distributor	driver
dementia	diana death	district	driver seat
democrat	diana fatal	disturb	drop
demograph	diarrhea	divid	drug

drug abus	editori	energi	eunic simpson
drug administr	editori accompani	enforc	europ
drug keep	educ	engag	european
drug keep angioplasti	educ need	engin	evad
drug keep angioplasti arteri	edward	england	evalu
drug keep angioplasti arteri clear	effect	england journal	even
drug regimen	effort	english	event
drug resist	egg	enhanc	eventu
drug therapi	eighth	enjoi	evid
dry	eighth season	enjoyabl	evid suggest
duchess	elderli	enrol	evolutionari
durango	elec	ensembl	evolutionari driver
durango need	electr	enter	evolutionari driver seat
durango need work	electron	entertain industri	evolv
durban	electron commerc	entertain present	exactli
dutch	electron edition	entertain report	exam
dutch studi	elem	entertain tonight	examin
dying	elev	enthusiast	examin med
dynamic	elimin	entir	examin med school
ear	elizabeth	entri	examin med school admis
earli	ellei	environ	examin med school admis criteria
earli onset	ellen	enzym	exampl
earlier	emerg	epic	excell
earlier relat	emerg contracep	epidem	excess
earlier relat stori	emi	epidemiologi	exchang
earn	emmi	episod	exclus
easi	emotion	equal	exec
easier	empir	equip	exec produc
easili	emploi	equiti	execut
east	enabl	era	execut produc
eat	enceph	erasmu	exerc
econom	enceph alert	erasmu univers	exercis
economi	enceph alert broaden	eric	exist
economist	encephalopathi	erickson	expan
edg	encount	especi	expand
edgar	encourag	essenti	expect
edinburgh	end	establish	expenditur
edit	endear	estat	expens
edition	endear perform	estim	experi
editor	endometri	estrogen	experienc
	endometri cancer	eunic	

experiment	fast track	fil	follow
expert	fast track trad	film	food
expert conclud	fat	film academi	football
expert contend	fatal	film critic	forc
expert sai	fatal human	film director	ford
expert warn	fatal human disord	film fest	ford motor
explain	father	film festiv open	forecast
explan	favor	film produc	foreign
explor	fax	filmmak	foreign box
expos	fda	fin	foreign box offic
exposur	fda overhaul	final	foreign total
expres	fear	financ	form
express	fear stir	financi	formal
extend	featur	find	forrest
extent	featur credit	find suggest	forth
extra	featur film	fingerprint	fortun
extract	feb	finish	forum
extrem	feder	finish stir	forward
eye	fee	finish stir melt	foster
fac	feel	finish stir melt pot	found
fact	feet	finland	foundat
factor	fell	firm	founder
factor ti	fellow	fit	fourth
fai	felt	fla	fourth annual
fail	femal	flar	fourth quarter
failur	fergi	flat	fox
fair	ferguson	fled	fractur
fall	fernandez	flem	franc
falter	fertil	flesh	francisco
fam	fest	flexibl	frank
famili	festiv	fli	frasier
famili histori	festiv open	florida	fred
famili member	fet	flower	french
famili orient	fewer	flu	frequent
familiar	fiction	flu epidem	fresh
famou	field	flu shot	fridai
fan	fifth	flu shot urg	fridai octob
fantasi	fifth annual	focu	friend
far	fight	focus	friend wed
fashion	fight cancer	fold	fring
fast	figur	folk	front

fruit	gener motor corp	got allergi try	guid
fulli	genet	got allergi try wash	guid seri
fulvio	genet engin	gotham	guidelin
fulvio cecer	genet mutat	gotten	guild
fun	genr	govern	guilti
function	genuin	grab	guitarist
fund	georg	grad	gun
fund rais	georgia	gradual	gynecolog
funer	geriatr	graff	gynecologi
funni	germani	graham	habit
futur	get	grai	hair
futur flu	gianni	grai school	hair loss
futur flu epidem	gianni versac	gram	hair loss possibl
gai	giant	grand	half
gai men	giant said	grand slam	half hour
gain	gift	grant	hall
gala	girl	graphic	halt
galleri	girlfriend	great	hand
gam	giv	greater	hand heart
gang	given	greater risk	hand heart healthi
gang relat	global	greatest	hander
gang relat mark	global media	greatli	handl
gang relat mark shakur	go	greek	hang
gap	goal	green	happen
garcia	goe	greg	happi
gari	gold	gregori	hard
gari graff	goldberg	grew	harder
gartner	gon	gross	harm
gartner group	gonna	ground	harpercollin
gastrointestin	gonna plai	group	harri
gat	gonorrhea	group continu	hart
gather	gonorrhea infec	grow	harvard
gav	gonorrhea infec increas	grown	harvard medic
gen	good	growth	harvei
gen expres	good new	guarante	hat
gen link	goodman	guarante impli	haunt
gender	gor	guarante impli regard	hav
gener	gordon	guarante impli regard inclusion	haye
gener believ	gorman	guard	hazard
gener hospit	got	guest	head
gener motor	got allergi	gui	heal

health benefit	henri	hollywood film	hyperten
health car	henri paul	hollywood film festiv	hypothes
health car burden	hepat	holocaust	ibm
health car cost	hercul	hom	ic
health car worker	heroin	homer	ic storm
healthi	herp	honor	icon
healthi stat	herpesviru	honore	ida
healthi stat list	herpesviru link	hontz	idea
healthier	herv	hop	ideal
hear	heston	hopkin	ident
hear aid	hid	hopkin univers	identifi
heart	high	hormon	ignor
heart attack	high level	horror	iii
heart diseas	high rat	hors	ill
heart failur	high risk	hospit	illeg
heart healthi	high school	host	illinoi
heart muscl	high speed	hot	imag
heart muscl cell	higher	hotel	imagin
heart patient	highest	hotlin	immedi
heart transplant	highest rat	hour	immers
heart transplant rejec	highli	hous	immun
heat	highlight	hous panel	immun cell
heavi	highlight octob	hous pass	immun practic
heavier	hill	hous pass fda	impact
heavili	hind	hous pass fda overhaul	impair
height	hir	household	implant
held	histori	household averag	impli
helen	hit	household rat	impli regard
hell	hitchcock	hug	impli regard inclusion
helm	hitter	hugh	implic
helmer	hiv	human	imporov
help	hiv infec	human disord	import
help crohn	hiv posit	human servic	import produc
help crohn diseas	hiv posit worker	humor	impress
help hand	hmo	hundr	improv
help hand heart	hmo perform	hunt	in
help hand heart healthi	hmo perform vari	hurrican	inaugur
help parkinson	hmo perform vari wid	hurrican fear	incid
help prevent	hold	hurrican fear stir	includ
help prevent airbag	holidai	hurt	inclusion
help prevent airbag injuri	hollywood	husband	incom

increas	inser	interven	jennif
increas risk	insert	interview	jeopardi
increasingli	insid	intl	jerri
incub	insid daili	intric	jersei
incub period	insid daili varieti	intrigu	jesu
independ	insid daili varieti highlight	introduc	jewish
independ featur	insid daili varieti highlight	invest	jim
indi	octob	investig	jimmi
indian	insid edition	investig sai	joan
indic	insight	investig target	job
individu	insomnia	investor	job said
induc	insomnia rais	involv	joe
industr	insomnia rais health	involv patient	joe pesci
industri	insomnia rais health car	ion	joel
ineffect	insomnia rais health car burden	ion rat	johannesburg
infant	inspir	iowa	john
infant liv	instanc	iron	john hopkin
infant stor	institut	irradi	john hopkin univers
infec	instrum	isol	john wort
infec increas	insur	issu	johnni
infect	intak	itali	johnson
infect patient	intak prolong	italian	joi
infecti	intak prolong lif	item	join
infecti agent	intel	jack	jon
inflamm	intens	jack nicholson	joseph
inflammatori	interact	jackson	journal
inflat	interest	jacqu	journal natur
influenc	interim	jai	journal pediatr
influenza	interim chief	jail	journal scienc
inform	interim chief execut	jakob	journalist
infrastructur	intern	jakob diseas	judg
infusion	intern documentari	jam	judg rul
inherit	intern film	jama	juli
inhibitor	intern film festiv	jan	julio
initi	intern film festiv open	januari	julio martinez
injec	intern medicin	japanes	jump
inject	intern sal	jazz	jun
injur	internet	jean	jungl
injuri	internet access	jeff	jurass
ink	internet commerc	jenni	jurass park
inocul	internet sit	jenni hontz	juri

juri order	ladd	leagu championship seri	literatur
juri order chrysler	ladi	learn	literatur priz
just	lak	learn word	littl
just wrap	lak buena	leav	littleton
justic	lak buena vista	led	liv
justifi	lamb	lee	lloyd
kathi	land	left	lloyd webber
kathi lee	lang	leg	lo
kathleen	languag	legal	lo angel counti
keep	lap	legisl	lo angel lo
keep angioplasti	larg	lend	lo angel lo angel
keep angioplasti arteri	larger	length	lo angel trial
keep angioplasti arteri clear	largest	lens	local
kei	larri	leonard	locat
kei adult	last	leonard kladi	lolita
ken	lat	lesli	lon
kennedi	latest	lesli lang	london
kept	latest movi	lesson	london british
kevin	latex	lethal	long
kevin smith	latex label	letter	long liv
kick	latex label requir	level	long term
kid	latin	levi	long tim
kil	latino	lewi	longer
kill	laugh	li	longer period
killer	launch	liar	look
kind	laura	liar liar	loom
king	lauren	lif	lopez
kiss	law	lif threaten	los
kiss plant	lawmak	lifestyl	loss
kladi	lawrenc	lifetim	loss possibl
know	lawsuit	lifetim achiev	lost
knowledg	lawyer	lifetim achiev award	lost world
known	lax	light	lot
korea	lead	likelihood	loui
korea second	lead author	lin	lov
korea second seoul	lead caus	linda	lover
label	lead studi	lineup	low
label requir	lead studi author	link	low calor
labor	leader	lisa	low calor intak
laboratori	leagu	list	low calor intak prolong
lack	leagu championship	listen	low calor intak prolong lif

lower	mark	medic associ	mic genet
lower breast	mark shakur	medic center	michael
lower breast cancer	marker	medic cost	michael flem
lower breast cancer risk	market	medic devic	michael jackson
lukewarm	marri	medic journal	michell
lukewarm profit	martha	medic record	mid
luncheon	martin	medic school	middl
lung	martinez	medic treatment	middl ag
luxuri	marv	medicin	midst
lynn	marv albert	meet	migrat
lyric	marvel	melani	mik
machin	maryland	melod	milan
mad	mass	melt	mild
mad cow	massachusett	melt pot	militari
mad cow link	massachusett gener	member	miller
madison	massachusett gener hospit	membership	mind
magazin	match	memor	minim
magic	materi	memori	minnesota
mail	matern	men	minnesota top
main	matter	men behav	minnesota top healthi
mainli	matthew	men behav badli	minnesota top healthi stat
maintain	maximum	menac	minnesota top healthi stat list
majendi	mcbeal	menopaus	minor
major	mccarthy	menstrual	minut
mak	mccartnei	menstrual period	mirag
maker	mean	mental	miramax
mal	measur	mental health	miramax film
malibu	meat	mention	misdemeanor
mall	meat product	mer	misdemeanor assault
malnutri	meat recall	merced	miss
man	meat recall run	merced benz	miss diana
manag	mechan	merchand	miss link
manhattan	mechan underli	messag	mississippi
manner	med	met	mistakenli
manslaught	med school	metabol	mittchell
manufactur	med school admis	method	mix
march	med school admis criteria	methodist	mix set
margin	media	mexican	model
mari	media access	mexico	moder
maria	media access award	miami	modif
marin	medic	mic	modifi

mold	movi	near	nobel literatur priz
molecul	movi theater	nearli	nobel priz
molecular	multipl	nearli ident	nomin
molli	multipl sclerosi	nearli undetect	non
mom	mundell	necessari	nonsmok
moment	murder	necessarili	normal
mon	murdoch	necessarili mean	north
mon oct	murdoch said	neck	north carolina
mon sep	muscl	need	northern
mondai	muscl cell	need work	not
mondai octob	muscl tissu	neg	notabl
mondai septemb	museum	negoti	notabl quot
monei	music	neil	notic
monica	musichound	nelson	nov
monica roman	musician	neonat	novel
monitor	mutant	nerv	novemb
month	mutant gen	nerv cell	nsaid
monti	mutant gen link	netherland	nsaid user
mood	mutat	network	number
moon	mycobacterium	neurolog	nurs
moor	mysteri	neurologi	nutrition
moral	mysteriou	neuron	nypd
morbid	mysteriou second	new corp	nypd blue
morn	nab	new coverag	obes
morri	nam	newborn	obes risk
mortal	nanci	newcom	oblig
mortal weekli	narr	newli	obscur
mortal weekli report	narrat	newsmagazin	observ
morton	nasdaq	newspap	obstetr
mother	nation	nic	obtain
motion	nation academi	nicholson	occasion
motion pictur	nation averag	nick	occup
motiv	nation health	nicol	occur
motor	nation institut	nielsen	oct
motor corp	nation syndicat	night	octob
motor function	nation syndicat journalist	nih	octob issu
motown	natur	nih studi	octob wed
mountain	natur medicin	ninth	octob wed oct
mountain view	nazi	noah	odd
mov	nba	nobel	odonnell
movem	nbc	nobel literatur	offer

offic	organ	parkinson diseas	perfect
offic said	orient	parkinson patient	perform
offici	origin	part	perform vari
offici sai	oriol	parti	perform vari wid
offspr	oscar	partial	period
ogea	out	particiap	permiss
old	outbreak	particip	permiss parent
old ag	outcom	particl	persist
older	output	particular	person
older american	output start	particularli	person comput
older peopl	outrag	partner	pertussi
oldest	outsid	partnership	pesci
oldest old	overhaul	pasadena	peter
oliv	oversea	pass	petrikin
oliv ston	overwhelm	pass fda	philip
on	ow	pass fda overhaul	phillip
onair	own	passion	phon
oneill	owner	past	photo
ongo	oxid	past year	photograph
onlin	pac	pathogen	physical
onset	pacif	patient	physician
ontario	pack	patrick	physiologi
onward	pact	pattern	phytoestrogen
onward christian	pai	paul	phytoestrogen lower
onward christian soldier	paid	paul majendi	phytoestrogen lower breast
open	pain	paul oneill	phytoestrogen lower breast
open sept	paint	peacemak	cancer
oper	pair	pediatr	phytoestrogen lower breast
oper system	palo	pediatrician	cancer risk
opera	palo alto	pen	piano
opinion	panel	pennsylvania	pic
opinion express	paparazzi	pensacola	pick
opportun	paparazzi photograph	pentium	pictur
oppos	paper	peopl	piec
opposit	paper said	percent	pig
oprah	parallel	percent profit	pill
oprah winfrei	parent	percentag	pioneer
option	pari	perdita	piston
orchestra	park	perdita durango	pitch
order	parker	perdita durango need	pix
order chrysler	parkinson	perdita durango need work	plac

placebo	posit patient	presenc	priz
placem	posit worker	present	prob
plai	position	preserv	probabl
plan	possess	presid	problem
plant	possibl	presid clinton	procedur
plastic	possibli	press	proceed
playboi	post	press confer	process
player	post lukewarm	press releas	prod
playoff	post lukewarm profit	pressur	produc
playoff gam	postseason	preval	produc compani
plaza	pot	preven	produc david
plc	potent	prevent	produc design
plead	potenti	prevent airbag	produc jam
pleasur	potenti fatal	prevent airbag injuri	product
plenti	poulti	prevent car	profession
plot	pound	prevent malnutri	professor
plu	poverti	preview	profit
plung	power	previou	profound
poet	practic	previou stori	program
poetri	prais	previous	program content
poetri dai	pre	pric	progres
point	predict	primari	progress
point repres	predict heart	primari goal	project
polanski	predict heart transplant	primetim	prolong
polic	predict heart transplant rejec	primetim liv	prolong lif
polici	predictor	princ	promis
polish	prefer	princess	promot
politician	pregnanc	princess diana	promotion
poor	pregnanc diabet	princip	prompt
poorli	pregnanc diabet risk	principl	properti
pop	pregnanc diabet risk identifi	print	propos
pop star	pregnanc factor	prion	prosecutor
pop star michael	pregnanc factor ti	prion diseas	prosecutor offic
pop star michael jackson	pregnanc weight	prion research	prospect
popul	pregnanc weight gain	prion research win	protec
popular	pregnant	prion research win nobel	protect
portion	prematu	prion research win nobel priz	protein
portrait	premier	prior	protein ti
portrayal	premier week	prison	prov
pos	prepar	privat	proven
posit	prescrip	privat lif	provid

provid import	rapidli expand	regul	request
provinc	rar	regular	requir
psychiatri	rat	regular season	requir medic
psycholog	rat point	regularli	rerun
psychologi	rat point repres	rejec	rescu
public	ration	reject	research
public health	reach	rel	research admit
publish	react	relaps	research believ
purchas	reaction	relat	research conclud
purpos	read	relat death	research contend
pursu	readi	relat deficit	research institut
push	real	relat diseas	research led
put	realli	relat ill	research not
qualiti	reason	relat mark	research point
quarter	rebound	relat mark shakur	research report
quarter earn	recal	relat stori	research sai
quartet	recall	relationship	research sought
queen	recall run	relax	research stat
quest	receiv	releas	research win
questionnair	receptor	reliabl	research win nobel
quickli	recipi	relief	research win nobel priz
quiet	recogn	religi	research writ
quirki	recommend	remain	resembl
quit	record	remark	resid
quot	recov	rememb	resist
rac	recoveri	remind	reson
radic	red	remov	resourc
radio	red meat	render	respect
rai	reduc	rent	respond
rai richmond	reflect	rent collect	respons
rais	refus	rent collect prais	rest
rais health	regain	rep	restric
rais health car	regard	repair	restric slow
rais health car burden	regard inclusion	repeat	restrict
randomli	regardless	replac	restrict diet
randomli assign	regi	replic	result
rang	regi kathi	report	retail
ranger	regi kathi lee	report sai	retain
rank	regimen	reportedli	retard
rap	regina	repres	retir
rapidli	region	republican	return

reveal	ron	salem	scrib
revenu	room	salli	script
revenu com	ros	sampl	scripter
revers	rosenberg	san	se
review	rosi	san francisco	sean
review articl	rosi odonnell	santa	season
review oct	roster	santa barbara	season premier
review sept	rotterdam	santa monica	seat
reynold	round	sarah	seat risk
rhythm	round pick	sarah ferguson	seat risk infant
rich	routin	satellit	seat risk infant liv
richard	rowan	satir	seatbelt
richmond	rowan atkinson	satisfactori	seatbelt help
rid	rowan atkinson comedi	satisfi	seatbelt help prevent
right	royal	saturdai	seatbelt help prevent airbag
right hander	rug	saturdai night	seatbelt help prevent airbag
rigor	rul	sav	injuri
ris	run	saw	seattl
risk	rupert	scandal	seattl marin
risk factor	rupert murdoch	scen	second
risk identifi	rusedski	scenario	second album
risk infant	rusedski serv	schedul	second annual
risk infant liv	ryan	school	second dos
ritz	sad	school admis	second quarter
rival	saf	school admis criteria	second round
road	safer	scienc	second seoul
rob	safeti	scientif	second tim
robert	safeti expert	scientif evid	second week
rock	safeti rul	scientist	secret
rock band	sai	sclerosi	secretari
rock star	sai research	scop	secretari donna
rodent	said	scor	secur
roger	said mondai	scotland	seduc
rol	said studi	scott	seed
rol ston	said thursdai	scott erickson	seek
roll	said tuesdai	scream	seen
roman	said wednesdai	screen	seen post
roman polanski	saint	screen prevent	seen post lukewarm
romanc	sal	screen prevent malnutri	seen post lukewarm profit
romant	sal law	screenplai	seg
romant comedi	salari	screenwrit	seinfeld

select	sham	sir paul	softwar
self	shampoo	sister	softwar compani
self report	shap	sit	softwar develop
sell	shar	sitcom	softwar giant
seller	sharon	situat	sol
seminar	sharp	sixteen	sold
sen	shed	sixth	soldier
senat	shelter	siz	solid
send	shin	skill	solo
senior	shirt	skin	solution
senior columnist	shoot	skin reaction	somewhat
senior investig	shop	sky	son
sens	short	slain	song
sensibl	shot	slam	soni
sensit	shot urg	slat	soon
sent	shoulder	sleep	sort
sentenc	show	sleepless	sought
seoul	showbiz	sleepless children	soul
sep	showbiz event	slightli	soul food
separ	showcas	slot	soul man
sept	shown	slow	sound
septemb	sick	slow ag	soundtrack
sequel	sid	slow ag relat	soup
seri	sight	slow ag relat deficit	sourc
seri open	sign	slow parkinson	sourc clos
serious	signal	slowli	sourc said
serv	signific	small	south
server	significantli	smaller	south africa
servic	silenc	smart	south african
servic offer	silicon	smith	south carolina
servic provid	silver	smok	south korea
session	similar	smok cessat	south korea second
set	simon	smoker	south korea second seoul
set decor	simpl	smoker quit	southern
set design	simpli	snappi	spain
seventh	simpson	soap	span
sever	sinai	soar	spanish
sex	sing	social	spar
sexual	singer	societi	spark
shakur	singl	soderbergh	spawn
shalala	sir	soft	speak

speci	star	straight	subscrip
special	star michael	strain	subsequ
specif	star michael jackson	strand	substanc
spectrum	star trek	strap	substanc known
specul	starrer	strateg	succe
speech	start	strategi	success
speed	starter	streak	sudden
spend	stat	street	suffer
spent	stat list	street journal	sugar
spielberg	statem	strength	suggest
spin	statem releas	stress	suicid
spinal	statem said	stretch	suit
spinal cord	station	strik	summari
spirit	statist	strip	summer
spok	statu	stripper	sun
spokesman	steadi	strok	sun said
spokesman said	step	strong	sun system
spokesperson	stephan	strongli	sundai
spokeswoman	stephen	structur	sundai night
spongiform	steril	struggl	sundanc
spongiform encephalopathi	steroid	student	sunset
sponsor	stev	studi	super
sporad	stev gorman	studi appear	superior
sport psychologi	stev jam	studi author	superman
sportscast	stev job	studi author conclud	supernatur
sportscast marv	steven	studi confirm	supervis
sportscast marv albert	steven spielberg	studi involv	supervis editor
spot	stewart	studi involv patient	supervisor
spread	stimul	studi lead	supplem
spread tuberculosi	stir	studi lead author	suppli
spring	stir melt	studi particip	support
squar	stir melt pot	studi provid	suprem
stabl	stock	studi show	suprem court
staff	stockholm	studi suggest	sur
stag	stomach	studio	surfac
stai	ston	stuff	surgeon
stak	stop	stunt	surgeri
stallon	stop smok	style	surgeri help
stand	stor	subject	surgeri help parkinson
stand ston	stori bas	submis	surgic
standard	storm	submit	surgic remov

surpris	technicolor	theresa tamkin	tom
surprisingly	techniqu	thermostat	ton
surround	technolog	thesp	toni
survei	technolog	thing	toni scott
surveill	technolog	think	tonight
surviv	teen	third	took
surviv colicki	teen smok	thoma	took plac
surviv colicki babi	telecommun	thought	tool
survivor	telephon	threat	toothless
suscept	televi	threaten	top
suspect	televi art	thriller	top healthi
sustain	televi art scienc	thriv	top healthi stat
sweden	televi market	thursdai	top healthi stat list
swing	televi	thursdai night	toronto
switch	televi ceremoni	thursdai octob	toronto film
sylvester	tell	thursdai septemb	toronto film festiv
sylvester stallon	tend	ti	total
sympathet	tennesse	ticket	touch
symphoni	teresa	ticket sal	tough
symposium	term	tightli	tour
symposium host	termin	till	tournam
symptom	terrif	tim	town
syndicat	territori	tim greater	track
syndicat journalist	test	tim greater risk	track trad
system	test reveal	timeslot	tract
tak	testicular	timi	trad
taken	testicular cell	tip	tradition
tal	testicular cell slow	tissu	tragedi
talent	testicular cell slow parkinson	titan	tragic
talk	testimoni	titl	train
taller	tetanu	titl charact	transac
tamkin	tetanu pertussi	tml	transform
tap	texa	tobacco	translat
target	texa ranger	todai	transmis
target shelter	theater	todd	transmit
task	theatr	todd mccarthy	transplant
tax	theatric	told	transplant patient
taylor	theoret	told daili	transplant rejec
team	theori	told daili varieti	treat
tear	therapi	told report	treat addic
technic	theresa	told reuter	treat hyperten

treat patient	typical	user	villag
treatment	ulcer	usual	violat
trek	ultim	uterin	violenc
trend	unabl	vaccin	violent
tri	unbridl	valu	violent commerci
trial	unclear	valuabl	violent commerci evad
tribut	uncontroll	van	viral
trigger	undergo	vancouv	virgin
trio	underli	vari	virginia
trio finish	understand	vari wid	virologi
trio finish stir	understood	varieti	virtual
trio finish stir melt	underw	varieti chief	viru
trio finish stir melt pot	undetect	varieti chief film	virus
troi	unexpected	varieti chief film critic	vision
troi woefully	unfocus	varieti cyber	visit
troi woefully unfocus	unfortun	varieti cyber summari	vista
troubl	union	varieti entertain	visual
troup	uniqu	varieti entertain report	vital
true	unit	varieti highlight	vitamin
true stori	unit stat	varieti highlight octob	vocal
truli	univers	varieti senior	vocalist
truth	univers colleg	varieti senior columnist	voic
try	univers medic	variou	volum
try wash	univers pictur	vast	vot
trying	univers school	veget	wai
tub	unix	vehicl	wait
tuberculosi	unknown	venou	wak
tue	unnecessari	ventur	wal
tue oct	unpredict	venu	walk
tue sep	unprotect	versac	walker
tuesdai	unusu	version	wall
tuesdai octob	up	versu	wall street
tuesdai septemb	upcom	veteran	wall street journal
tumbl	upcom art	vic	walter
tumor	upcom film	victim	want
tun	updat	victor	war
tunnel	upn	victori	war iii
turn	upset	video	ward
twic	urban	video confer	warn
twist	urg	view	warner
type	us	viewer	warner bro

wash	went	wir	worsen
washington	west	wis	wort
watch	western	wisconsin	worth
water	westport	wit	wrap
we	wheel	woefulli	writer
weak	whit	woefulli unfocus	written
wear	whit hous	woman	wrong
wear hear	whoopi	women	wrot
wear hear aid	whoopi goldberg	women report	xofficeoversea
weather	wid	wonder	yanke
webber	widescreen	wonder world	yard
wed	wif	wor	year
wed oct	wil	word	year ago
wednesday	wild	work	year ago averag
wednesday octob	william	worker	year contract
wednesday septemb	wilson	workplac	year period
week	win	workstat	year studi
week issu	win nobel	world largest	york
weekend	win nobel priz	world premier	york citi
weekend box	wind	world seri	york film
weekend box offic	window	world trad	york stock
weekli	winfrei	world war	york yanke
weekli report	wing	world war iii	young
weigh	winner	world wid	young peopl
weight	winston	worldwid	younger
weight gain	winston salem	worn	zero
welcom	winter	worri	

ข.4 พืเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลีของชุดเอกสาร Reuters-Top10

พืเจอร์ที่ได้จากการเลือกโดยใช้ต้นไม้วลีจากชุดการทดลองทั้ง 4 ชุด ได้จำนวนพืเจอร์ 3,421 พืเจอร์ ดังนี้

abandon	abil	access	accumul export
abbott	abl	accord	achiev
abbott laboratori	abruptli	accordingli	acknowledg respons
abegglen said	absolut sens	account	acquir
abid	accept	accru	acquir total explor australia

acquisi	aggress	alloc	announc
acreag	aggress	alloc rat	annual
acreag bas subject	aggress seek acquisi	allow	annual level meet
acreag reduc	ago	allow bank	annual meet
act	agoni	allwast	annual sal
act quickli	agre	altern	annual sharehold
action	agre trad deal	aluminum	answer
activ	agreem	alver carlson	anticip
activ central	agreem follow month	alw	appeal
actual damag	agreem guarante access	amend	appear
ad	agreem sign	america	appli
ad tom campbell	agreem total	america attack	applic
add	agricultur	america said	appoint
addition	agricultur depart	america simpli react	approach america
addition mln dlr	agricultur depart gav	american	appreci
addition quota	agricultur offici	american citizen liv	approach
address excess consump	agricultur produc	american dollar	appropri
adelaid	agricultur report said	american farmer	approv
adher	agricultur secretari	american forc	apr argentin grain board
adjust	agricultur secretari carlo	american helicopt cam	preliminari
adjust export pric	dominguez said	american liv abroad	april
administr	ahead	american peopl	april cheap oil feedstock
administr offici	aid	anaheim	april european currenc market
admit	aim	analyst	react quietli
adop	air	analyst ad	april japan littl known ministri
advanc	alan wheatlei	analyst believ mondai wall	april kuwait
advers affect	alaska	street crash	april pest
advers effect	alaska colvill delta	analyst darrell	april qtr
advic	alaskan oil trad continu	analyst expect	april remark
adviz	albeit	analyst expect bonn	april stood
adviz	alberta	analyst rang	april taiwan import
advoc min	ald	analyst sai	april todai turmoil
aerospac	alert	analyst said	april union texa petroleum said
affair	algeria	analyst said west germani	april year
affair expert	ali	stubborn march	arabia
affect	alleg	anbaa	arbitr
affili	alli	andov	area
africa	alli accord	anger	argentin
africa nil nil nil nil	alli depend	angri	argentin grain board adjust
aftermath	alli leader	anhui	argentina
agenc	alli signal	anniversari	argentina design

argentina said	australia	bank interven	big
argu	australian	bank meet	biggest
ark	australian labor union	bank presid dieter hiss	biggest reduc
arm embargo	australian port	bank said	bill
armen estim	australian westpac bank	banker	block
arrang	australian wheat export	banker said	blockad
ash	author	bankruptci	blockad iran
ashland oil	author common shar	bar	blown
ashraf fouad	authoris	bargain	board
asia	automat pilot	barlei	board approv
asid	automobil manufactur	barrel	board declar
ask	avail	bas	boat
ask congression leader	avail opportun	basi	bolder step
ask gatt	aveng	basic commod petrochem	bolster
ask presid	averag	battl	bomb threat
assess	averag daili reserv	bavadra	bomb threat forc evacu
asset	avg shr	bbl	bond
assist	avoid	beacon	bond market
assist act	awai foreign investor	beacon usual lit	bond suffer
associ	awar	bear	bonn
associ profit	backdrop	bearish	bonn refus
assum	bader told	beef	bonu issu
assump	badli	beer	boost
atlant	bahrain	began	boost output
atlanta	baker	begin	boost rat
atpp block	baker appear	begun escort	bor
attack	baker cut short	behalf	borg warner
attack involv	baker said	believ	borg warner corp
attack kuwaiti connect vessel	baker statu	believ econom	borrow
attack respond	baker weekend blast	bell	boston
attempt	baker weekend remark	benchmark crud west texa	boston corp
atten	bal	intermedi	bottom
attend	balanc	benefit	bought
attitud	ball rol	bentsen said	bpd
attract	ban	best	bracket
attribut	ban impos	best servic	bran pollard wheat
aubrei	bangkok	best year	brav
audienc	bank	better	brazil
aug	bank borrow	beverli hill	brazil suspend payment
august	bank economist	bicol	brazilian loan
australasia	bank economist said	bid	breach

bread wheat	bundesbank presid karl otto	cash	chemic
bread wheat nil nil nil	pochl	cash crop	chemic acquisi
break	burden	casualti	chemic busi
brief	bureaucrat	categori paper todai	chemic engin
bring	bureaucrat counterpart	caught lai min	chemic execut shar
bring suppli	burlington	caus	chemic export
britain	burlington industri	caution	chemic industri
britain cabl	burnham lambert	cdi	chemic industri biggest custom
british	bushel	ceil	chemic maker
british chancellor	busi	celebr	chemic manufactur
british petroleum	businessmen	cellular	chicago
british print	buyer	cent	chicago manufactur
broadli	cain said	center	chief
broker	cal	centr	chief economist
brother	calcul	central	chief execut offic
brought	calgari	central bank	china
brunt	calif	central bank interven	china daili said
brussel	california	central banker	china nil nil nil nil
buck	call	central banker rais	china wheat crop threaten
budget	calm	central oklahoma	chines
budget deficit	cam	cereal	chip
budget deficit rel	cambridg	cereal substitut	choic
budget mln dlr	campaign	certain	chok
budget sav propos draft	campeau	certain fish	chokepoint
bueno air	canada	certainli	christma tree
bui	canada februari trad surplu	chairman	christoph hanson
build	canadian claim	chairman gordon cain	cincinnati
build block	canadian dlr	chairman patrick leahi	circumst
builder	cancell	challeng	cit
built	capabl	chancellor	citi
bullish	capit	chang	citicorp
bullish outlook	capit spend	character	citru fruit
bullish view	capit stock	charg	civil defens offici said
bundesbank	car	charg off	civil disobedi
bundesbank allot	career naval offic	charg off total	claim
bundesbank board member clau	cargo	charl	clash
koehler	carri	charter	class
bundesbank board member clau	carri awai	charter soviet tanker	clear
koehler yesterdai	carri kuwaiti oil	charter tanker	clearli
bundesbank latest liquid allot	carryforward	chas	clercq
bundesbank offici	cas	chas awai	clercq told

cleveland	compani	conflict	convert prefer shar
climb	compani common stock	conflict unduli	convert subordin debentur
clos	compani continu turn	confront	convinc
close link	compani produc	congress	cooper
closure	compani report	congression concern	cooper basin
coars grain	compani result	congressman	coordin interven
coars grain produc	compani said	conn	cor busi
coast	compani spokesman	connec	corn
coastal	compani tonnag	consensu	corn plantat
cohen	compani total energi produc	consequ	corn trader
collaps	compani went public	conserv	corp
colo	compar	consid	corp said
colombia	compar figur	consid neglig	corpor
com	compar quarter	consid propos	corpor purpos
combin	compar week	consider	correct
comment	comparison	consist	correspond
commerc	compens	consolid	cost
commerci	compet	consortia trying	cotton
commis	competi	consortium	cottonse oil
commis authoris	competit	construc	counti approv
commis cereal manag	competitor	consult	counti loan rat
commis sourc said	complaint	consult firm asia advisori	counti loan rat differenti
commission willi	complet	consum	countri
commit	complet acquisi	consumm	countri oust prim minist timoci
committe	complet chok	consumm politician	bavadra renew pressur
committe reject	complet merger	contain	countri sugar
committe spokesman told	completion	content	coupl
reuter	compli	context	cours
commod	complianc	continu	court
commod chemic	comprom	continu export	cover
commod chemic busi	comput	continu fight	cower
commod provid	concentr	continu talk	crack
commod report	concern	contract	crazi eddi
common	conces	contrast	creat
common secur interest shar	concret	contribut	credibl
common shar	condens	control	credibl strategi
common stock	condition	control congress	credit
common stock outstand	condition continu	control declin	credit suiss
commun	confer	conven	crisi
commun corp	confid	conver	critic
communiqu	confirm	convert	critic charg mpt

crop	custom design	decad	darman said
crop export registr	custom repurchas agreem	decemb	deregul
crud	cut	decid	describ
crud oil	cut abroad	decid factor	design
crud oil post	cut rat	decision	desir
crush blow	cycle	declar	despit
crush wheat	dai	declar destin	destabil drop
csi	dai account	declin	destin
csi said washington	dai averag	declin comment	destroi
ct	dai banker secur	declin earn	destroi kei bas
ct ct	dai drop	defenc	destroyer
ct loss	dai sama deposit	defens	destroyer kidd
ct loss ct	dai sama deposit ris	defens depart said just pct	destroyer pump
ct prior	dairi	defens secretari caspar	destruc
ct profit	dalla	weinberg told	detail
ct profit ct	damag	deferr	deter
cubic feet	damag hom	defici payment	deter iranian attack
cumul	dampen	deficit	deterior
cumul export	danger	deficit cut	determin
cumul figur export registr	dash	definit	deterr polici
cur	dat	definit agreem	develop
curb	dat interbank	delai	di
curb monei suppli growth	data	deleg	diagnost
currenc	david jon	deliv	differ
currenc analyst	dead	deliveri	differenti
currenc analyst refus	deadlin	demand	difficult
currenc declin	deal	democraci	dilemma
currenc market	deal agre	demonstr	dilig
currenc stabil	dealer	deni	dilut
currenc unit	dealer expect volum	denot minu figur	dilut shr
current	dealer said	denver	diminish
current discuss	dean witter reynold	depart	dip
current fiscal year	death	depart not	diplomat
current level	debat	depart offici said	diplomat mission
current manag	debentur	depart said	diplomat mission oversea
current oper	debentur holder	depart said produc	diplomat said
current rat	debilit trad war	depend	diplomat sourc
current season	debt	deposit	diplomat sourc said
current trad rang	debt burden	deposit account issu	direc
current year	debt burden dom	depreci	direct injec
custom	debtor countri export	deputi treasuri secretari richard	direct result

director	do	dupont conoco inc chemic busi	embark
director gener	dollar	durum wheat	embroil
disappoint	dollar collaps	earli	emerg
disast	dollar crisi	earlier	emerg loan
disclos	dollar declin	earlier clash	employe
disclosur	dollar-drop	earn	employe stock ownership plan
discontin	dollar fal	earn includ	employm
discontin oper	dollar fall	eas	enabl
discount	dollar fell	easili	enact
discount corp	dollar latest	east	encourag
discount corp mcgroarti describ	dollar open	east europ	end
discount rat	dollar open lower	east europ nil	end march
discount rat januari	dollar push	ec cereal manag	endang
discov	dollar slid	econom	energi
discoveri	dollar trend	econom crisi	engag
discreet rat adjust	dom	econom expan	engin
discus	dom petroleum	econom growth	england
discuss	domest	econom malais similar	enhanc
diseas	domest industri earn	econom polici	enjo
dismantl	domest market	econom polici decision	enorm
disput	domest polici	econom spotlight telecom	enquiri
distribu	dominion textil	econom stimul	enrol acreag
distributor	doshier said	econom summit	enroll
div	doubt	economi	entail world war
diversifi	dow	economi continu	enter
divest	dow chemic	economist	entertain
divid	dow jon industri averag	economist believ	entertain market
dividend	dramat cheaper	economist said	entir
division	drexel	ecu	entitl
dlr	driv	ecuador	entrench
dlr bid	driver seat	edelman	entri
dlr cash	drop	edmonton	equal
dlr compar	drought	edt	equip
dlr ct	drought account	effect	equiti
dlr gain	drug	effect jun	escal
dlr loss	dry	effort	escap
dlr loss ct	dry spell	elec	escort
dlr offer	dry spell continu	electron	escort kuwaiti
dlr profit	dump	elig	especi
dlr tax credit	dump semiconductor	elig winter wheat	establish
dmp	dupont	elimin	estim

ethylen	expd	facto	feed
eugen carroll	expect	factor	feedstock
eugen island block	expect econom stimulu packag	factor support	feel
euromarket	expect increas	fail	fell
europ	expell	failur	felt
europ nil	expens	fairli drawn	fertil
europ nil nil nil nil	expert told reuter tehran hold	fal	feudal
european	expert worri	fall	field
european commis	expir	fallen far	fifth year
european commun	expir midnight	famili	fight
european countri	explain	far	figur
european currenc unit	exploit	farina	figur includ
european market react quietli	explor	farm	fiji
european monetari	export	farm disast	fiji see
evacu	export commit	farm produc	fil
evalu	export commit usda	farm subsidi	final
evalu work	export hold option	farmer	final stag
eventu	export inspec	fat	financ
evid	export said	favor	financ minist
exacerb	export subsidi	favor technic chart signal	financ minist gerhard
exce	export tax	favour	financ offici
excess	express	fear	financi
exchang	expt	feasibl	financi advisor
exchang commis	extend	februari	financi condition
exchang loss	extens corpor restructur	februari april period	financi corp
exchang market	program	februari louvr accord	financi market
exchang rat	extent	fed	financi sector
exchang rat stabil	extern	fed expect	find domest polit prowess
exchequ	extern trad commission willi	fed fund	finish
exchequ nigel lawson	extra	fed interven	fir
exclud	extra chemic	feder	firm
exclud brazil	extra liquid	feder budget	firm said
exclud gain	extraordinari	feder disast relief	firm sal
execu	extraordinari item	feder republ	fiscal
execut	extrem	feder reserv	fiscal stimulu packag
execut offic	extrem danger	feder reserv bank	fiscal year
exist	fac	feder reserv chairman paul	fish
exlud	facil	volcker	fish disput
exp	fact	feder reserv enter	fit
expan	faction	feder sav	fitzwat said
expand	faction ignor mpt	fee	fitzwat told report

fix	fourth	gain addition trad right	govern offici
fix alloc rat	fourth quarter	gam	govern said
fix pric	fourth year	garo armen	govern secur market
fix rat	fragil	gatt	govern sourc said
fla	franc	gaug	govern welcom
flag	franc continu	gav	gra
flag kuwaiti tanker	franc said	gav iran	grab
flag kuwaiti tanker sea	francais de petrol	gener	grac
flag ship	frankfurt	gener agreem	grad
flag tanker	freedom	gener sav	gradual market absorp
flag tanker fleet	freefall	gener weak	grain
fledg war	freight rat	german	grain oilse registr
fleet	french	german bank	grand
fleet handl liquid cargo	french wheat	german monei	grant
float	fresh initi	german monei market split	great
flood	fresh veget	germani	greater
floor	friction	get	greater extent
florida	fridai	giant dow chemic	greek tragedi
flow	fridai averag	gibraltar	greenwich
flurri	friend	gillett	gross
focu	friendli	giordano said	ground
follow	frigat	giv	groundnutse oil
forc	frigat stark	given	group
forc attack	frontier	glamor	group led
forc baker hand	frontier insur	glenview	group said
forecast	frozen boneless beef	global econom	group shr
forecast australia	frozen hak fillet	global rat	grow
foreign	frtr	global recoveri mov	grow consensu
foreign currenc	frustrat	gnp	grow world
foreign exchang	fund	go	growth
foreign exchang market	fundament	gold	growth expect
foreign invest	futur	gold reserv	growth orient
foreign minist	futur belong	golden period	guard
foreign natur ga sal	ga	goldman sach	gulf
foreign purchas	ga field	gon ahead	gulf area
forgotten	ga properti	good	gulf follow
form	gaf	good prospect	gulf oil
fort worth	gaf chairman samuel heyman	gordon	gulf polici
fortun	estim	got	gulf reflag
forward	gaf heyman	govern	gulf war
forward look	gain	govern bond salesman	gulf war expert

hail	hemispher	hostil	import wheat
half	henan	hostil takeov attempt	impos
half quota	heng	hostil tender offer	impos mln dlr
halt	henlei	hour	imposi
hamper	henri bellmon	hous	impress
hand	henri schuler	hous speaker jim wright said	improv
happen	heublein	houston	improv inspec system
harbor	high	houston bas cain chemic	improv quota
harcourt	high alert	houston bas compani said	improv tariff
harcourt brac jovanovich	high margin busi	hug	improv trad condition
hard land	high qualiti beef	hug oil platform dot gulf	improv trad term
hard red spring wheat	high qualiti hilton beef	hurt	incent
hard red winter wheat	high stat	hyman	inch
hardli	higher	hyman said	inch shell
harlan ullman	higher cost	hyper inflat	includ
harri upham	higher pric	ibm	includ britain
hart scott rodino antitrust	higher rat	idaho	includ gain
improv act	hik	idea	includ loan loss provision
harvest	hill	identifi	includ meat
hastili cal meet	hint	identifi sal	includ sal
hav	hisanobu ohs	ignor	includ secur gain
hbj	histor experi	ill	incom
head	histori	ill advis	incom fund
headquart	hit	illeg	increas
health	hitherto	imbal	increas demand
heavi	hoax	imf	increas significantli
heavi import	hoel	immedi	increas tariff
hebei	hold	immedi impact	increas tension
hebei wheat aphid	hold compani	immedi shipment	ind
hectar	hold corp	imo	ind ttl
hectar valu	hold monei market tender today	impact	india
hectarag plant	holder	impair	india nil nil nil nil
hedg	hom	imperi oil	indiana
heighten	hong	implem	indic
held	hong kong	implement	indic total
held discus	honor	implicit	indic total includ report commit
helicopt destroi	honour	import	individu
helmet	honour pledg	import plac	indonesia
help	hop	import politician	industri
help block financ ministri plan	hormuz	import restric	industri analyst
help boost	horrif casualti	import shift	industri analyst said

industri nation	intern trad	iranian platform	jardin matheson
industri powerhous	internation	iranian ship	jiangsu
industri said	interven	iraq	jitteri investor
industri sourc said	interven just ahead	iraq blam iran	join
industri trend	interven period	iraqi missil	joint
industrywid	interven pric	ireland	joint statem
infantri attack result	interven rat	irna said	joint ventur
inflat	interven stor	ironi	jonathan lynn
inflationari	interven tender	irv bank	joseph
inflationari fear	interview	irv bank corp	journalist
inflationari fear boost gold	invad iran territori	isl citi	jpm
bullion	inventori	issu	judg
influcnc	invest	issu capit	juli
influenti polit group	invest firm	itali	juli vorman
inform	invest purpos	item	jun
initi	invest secur	jam	jun fiji today welcom
injuri	invest secur gain	jam abegglen	jun japan
insist	investig	jam oneill	jun onward
instabl	investor	jan	jun oversea shiphold group
institi	investor asher	januari	jun presid reagan said
institut	investor group	januari agreem	jun releas
insur	investor group led	januari provid	jun resid
insur compani	investor psychologi	japan	jun santo
integr	invit	japan car	jun year
inten	involv	japan economi	jun year end march
intend	iran	japan failur	just
intend acreag	iran iraq war	japan lucr oversea telephon	justic
intent	iran respons	busi	justifi
interest	iran rostam	japan main offer	kamal kharrazi
interest rang	iran sai	japan mov	kan
intermedi	iran start	japan nil	kansa citi
intern	iran war	japanes	keep
intern accord	iran war enemy iraq	japanes bank dealer said	kei
intern busi	iranian	japanes bank dealer warn	kei backer
intern econom	iranian attack	japanes electron	kei elem
intern forc	iranian gunboat	japanes export	kei industri sector
intern monetari gamesmanship	iranian man	japanes good	kei japan ministri
intern oil produc	iranian militari offici said	japanes insur	keiichi udagawa
intern preced exist	iranian new agenc irna quot	japanes market	kentucki bas oil refin
intern studi	iranian oil	japanes offer	kept
intern tanker	iranian oil platform	japanes offici	khorr fakkan

kil	lautoka	linse oil	lower rat
kil seamen	law	liquid	lower rat prospect
kill	law pass	list	lowest
kind	lawson	litig	lubbock yesterdai damag
kingdom	lawson said	littl	luck
knew	lawson told	liv	lucr industri
know	lawsuit	liv room	lukman
known	ldp	lo angel	luxembourg
koehler	lead	loan	machin corp
kohlberg kravi robert	lead rat	loan associ	macro econom
kokusai denshin denwa	leader	loan loss	magnitud
kong	leahi propos	loan loss provision	main
kotc	leas	local	main disput
kraft	leather footwear	local maiz import said	main showpiec
kuwait	leav	local offic	mainli
kuwait back	led	locat	maintain
kuwait oil tanker	left	london	mainten
kuwait tanker	leftwich	london televi interview tuesdai	maiz
kuwaiti newspaper	legal framework	long	major
kuwaiti oil export said	legal framework eas	long plan trip	major compani execut predict
kuwaiti oil tanker	legisl	long run stabil	major econom power
kuwaiti ship	lessen mpt author	long term	major impetu
kuwaiti tanker	letter	long term debt	major industri
kuwaiti tanker convoi	level	long term loan	major industri nation
lack	leverag buyout	long tim	major manag rol
lago	levi	longer	major oil
lai	liabil	longer term uncertainti	major plant capac
lanston	liber	look	major produc region
larg	liber cdi	loss	major shift
larg scal	libya	loss carryforward	major stak
larger	licenc	lost	major war
larger wheat set	licenc request	lot	mak
largest	lift	loui	maker
largest loss	light	louvr	manag
lat	likelihood	louvr accord	manag director
latest	limit	louvr agreem	manag director abdul fattah
latest attack seen point	limit access	louvr currenc accord	manhattan
latest econom packag	limit adjust	low	manila
latest exampl	limit retaliatori strik	lower	manufactur
latin american govern	lin	lower dollar	manufactur capac creat
launch	linda sieg	lower oil pric	march

march opec	mef	midnight	mln dlr budget
march show	meal	mil	mln dlr cash
margin	mean	mild	mln dlr loss
marin midland bank	measur	militari	mln dlr plu
marin midland oneill said	meat	militari action	mln dlr surplu
mark	medic profession liabil	militari coup	mln dlr worth
mark flow	medic quadrangl	militarili	mln hectar
mark rat	medicin	milwauke	mln loss
mark yesterdai	medium	min	mln mln
market	medium term	min iran	mln outstand
market believ	meet	min iran harbor	mln peso
market condition	meet need	miner	mln profit
market develop	melbourn	minimum	mln riyal issu
market fall	mellon	minimum export pric	mln shar
market instabl frighten	mellon bank chief economist	minist	mln shar outstand
market offer glimps	member	minist fail	mln stg
market pric	membership	minist kakuei tanaka	mln tonn
market psychologi	memori	ministri	mobil hom
market soft wheat	men	ministri action	moder
market suddenli believ fed	mention	ministri mov	modest
chairman paul volcker	merchant bank	ministri stanc	modest rat
market today	merg	minn	modest recoveri
marri	merger	minneapolis	modif
marri expect	merger agreem	minor	mondai
marri felt	merger propos	minor ministri	mondai attack
marri said	merrill	minu figur	mondai talk
marri told reuter	merrill lynch	minu total	mondai wall street
marshall aid plan	messeng	missil	monei
marshall plan	metr	missil strik	monei declin
martin	metropolitan plc	mission	monei mak petrochem busi
mass	mexico	mix	monei market
massiv	mich	mln	monei market dealer said
massiv budget deficit	michigan	mln acr	monei market oper
materi	mid	mln barrel	monei market rat
materiel	middl east	mln barrel report	monei market rat continu
matter	middl east expert sai	mln bpd	monei suppli
matter clos	middl eastern oil produc	mln canadian dlr	monetari author bought
matur	mideast	mln common shar	monetari polici
maxim	mideast gulf	mln debit	monitor
maximum	mideast gulf situat	mln dinar	monopol
maxwell	midland bank	mln dlr	monopoli

monsanto	naval escort	norman robertson	offer slightli higher rat
month	navi began	north america	offic
month ago	navi destroyer	north american	offici
month amid	navig	north chicago	offici figur
month declin	near	northern texa	offici iranian reaction
month includ	near record wheat crop	norwai	offici not
monticello	near term	not	offici pric
montreal	nearli	not compani	offici report
morgan	nearli halv	not includ bulgur	offici respons
morn	necessari	not total	offici sai
morocco	need	notabl	offici said
mortgag corp	negoti	notion	offici said minist
morton hyman said	negoti merger	novemb	offici washington sought
moscow	neogtiat tabl	number	offset
mothball ineffici plant	new	nurtur	ohio
mount	new confer	oaklei said	oil
mouth	newli elect	oat	oil compani
mov	newli releas annual report	object	oil dril platform
mpt	newspap	obstructionist	oil equival
mpt influenc	newspap quot	obtain	oil export
mtc	nich	obviou	oil field
mtc spent	nicosia	obvious	oil industri
mth	nigel lawson	occasion	oil industri analyst
mth includ gain	nigeria	occur	oil industri sourc said
mth includ loss	night	oceania	oil market
multilater	nightmar	oceania nil nil nil nil	oil minist
nadi	nikko secur intern	oct	oil output
nam	nil	oct allwast	oil platform
nap	nil nil	oct hug oil platform dot	oil pric
narrow	nil nil nil nil	oct iran	oil produc
narrow rang	non	oct kuwait oil export	oil rig mil east
nation	non accru loan	oct militari expert sai	oil stock ros
nation averag loan rat	non accrual	oct privat export report sal	oil tanker
nation econom summit	non accrual statu	oct rumor circul	oilse
nation largest independ oil	non expens	octob	oilse export pric adjust
nation total	non japanes oversea	octob shipment	oilse product
natsuo okada	telecommun firm	odd	okada
natur	non spanish barlei export	odd man	okla
natur disast	non tradition corn farm	offer	oklahoma citi
natur ga	normal	offer better trad term	oklahoma clean
natur lower	normal april	offer pric	old

old on	organ	par	peak
oljeselskap	origin	par valu	pek
opecc	origin bid	parent	pellet
opecc confer	orlando	pari	pend
opecc confer presid rilwanu lukman	osg see opec quota kei	pari accord	pension
opecc current pricc	ottawa	pari agreem	pension account
opecc meet	ounc	parti	pentagon
opecc member	outbreak	particip	pentagon announc
opecc produc	outgo	particular	peopl
opecc produc ceil	output	particularli	perelman
opecc produc quota drop	outrag	partli	perform
opecc ris	outstand	partner	period
open	outstand common stock	partnership	period end april
open market transac	outstand shar	pass	period includ
oper	outvot	pass legisl	permit
oper farm	overproduc	past	permit acreag
oper incom	oversea	pasta	permit atpp
oper loss	oversea own	patch	persian gulf
oper profit	oversea said pct	patent	pest
oper remain divid	overthrown	patrol gulf	petrochem industri
oper shr	overwhelm major dual	paul	petrochem plant
oper shr ct ct	own	paul oreffic	petrochem relat busi
oper shr loss	own liberian flag ship	payabl	petroleum
oper shr loss ct profit	own pct	payabl april	petromin
oper today	pac consult	payabl juli	pfennig
oper vessel	pacif	payabl jun	phas
opinion	packag	payment	philadelphia
opportun	pact	payout	philippin
oppos	pact expir	payroll	philippin coconut author said
opposit	pai	pct	coconut produc
option	pai april	pct annual	philippin damag agricultur crop
option origin sal plu actual	pai econom	pct annual rat	phoenix
export	paid	pct capac	pick
orchestr campaign	painstakingli slow declin	pct discount rat	piec
order	paint	pct duti	pipelin
ordinari	pakistan	pct gain	plac
ordinari shar	palai	pct leap	placem
oreffic	panic	pct ris	plai
oreffic said	paper	pct stak	plan
	paper said abnorm weather	pct stock dividend	plant
	condition	pct tariff	plant produc

plant util rat approach pct	pre tax profit	previous led	promin
plastic	precipit	pric	promin stem
platform	predict	pric cut	promis
plc	predict lower profit	pric increas	promot
pledg	prefer	pric tak	prompt gaf
plu	prefer dividend	prim	promptli
plung	prefer shar	primarili	propel
poehl	prefer stock	princip	properti
poehl not	prelud	principl	propos
point	premium	print	propos acquisi
point cut	prepar	prior	propos offer
pois	present	prior year	prorat
polic said	present level	privat	prospect
polic	present liquid monei market	privat center	protec
polic decision	condition	privat held	protect
polit	preserv	privat held cain chemic	protect gulf
polit analyst said	presid	privat held compani own	protect kuwaiti oil tanker
polit impact	presid reagan	probabl	protectionist postur
popular	presid reagan said today	problem	prov
port	presid reagan speak	proce	prov reserv fell
portfolio	press	proceed	proven
portug	press confer	process	provid
portug entri	press iraq	produc	provid consum
posit	pressur	produc bas	provinc
position	presum	produc loss	provision
positon hedg	pretax	produc restraint	psychologi turn
possibl	pretax earn	product	pty
possibl disturb	pretax profit	product area	public
possibl sal	prev	product pric	public adjust
possibli	prev week	profession	public offer
post	prevail	profit	publicli
postal	prevent	profit chemic	pull
postal sav	prevent massiv	profit specialti chemic	pump
postpon	preventedf	program	purchas
potent polit forc	previou	program crop	purchas author
potenti	previou week	progress	purchas pric
pound	previous	prohibit	purpos
powderi mildew	previous announc	project	pursu
power	previous announc acquisi	project corn	push
power polit faction	previous announc sal	project stock surpass	push profit margin
pre summit speech celebr	previous disclos	prolong dry spell	put

puzzl	rat fell	reduc	report outstand sal
qtly div	rat pul	reduc program	report quantiti
qtly div ct ct prior	reach	reduc trad	report said
qtr	react	reed	repres
qtr feb	reaction	reed shar	repres approxim pct
qtr mar	readi	refin	republican parti
qtr march	reaffirm	reflag	repurchas
qtr sept	reaffirm america willing	reflect	repurchas agreem
qtr shr	reaffirm support	reflect tighter suppli	repurchas pact
qtr shr ct ct	reagan	refus	repurchas pact rat soon
qualiti	reagan administr	regard	request
quantiti	reagan coupl	region	requir
quarter	reagan declar	registr	research
quarter earn	reagan hint	regret	reserv
quarter end	reagan impos	regul	reserv ad requir
quarter just end	reagan plan	regulatori	reserv bank
quarter oper	reagan polici	regulatori approv	reshadat oil platform
quarter profit	reagan said	regulatori control	resid return
quarterli dividend	real	reinforc	resign
question	real estat	reinstat	resist
question presid reagan	real estat asset	reiter advic	resolv
tomorrow	realiz	reject	respect
quickli	realli	reject licenc	respect market season
quito	reappoint	rekindl inflat	respond
quot	reason	rel	respons
quot today	recapit	relat	rest
quota	recapit plan	relat compani	restat
radio	receiv	relax	restaur
radio station	reces	releas	restor faith
rain	reciev	remain	restor investor confid
rainfall	reckon	remain shar	restric
rais	recogn	remark	restructur
rais cash	recommend	remedi	restructur program
rais fear	record	remov	restructur touch
rais hop	record april	renew pressur	result
rais rat	record dat	reorgan	result includ
rang	record year	repai	result reflect
rapese	recoveri	repeat past mistak	result restat
rat	red	repeat warn	resum
rat boost	redress gap world trad imbal	replac	retail
rat cut	redress global trad imbal	report	retain

retali	round	san	sector
retali continu	row	san diego	secur
retaliatori measur	royal gold	san mateo	secur repurchas
retir	rpt argentin grain	sanction	seek
retreat	rul	sank	seek acquisi
return	rul liber democrat parti	santo	seek control
reuter	rumor	santo group stak	seen
rev	rumour	santo said	seiz
revenu	run	satellit commun market	sel
revenu exclud	run bal	saturdai	sell
revers	sabah	saudi	semiconductor
revers stock split	saf	saudi arabia	semolina
review	sai	saudi arabian monetari	senat
revlon	said	sav	senat agricultur committe
revlon group	said approv	sav system abil	senat democrat
ric	said bach secur	saw	senat panel studi loan rat
richard	said dow	scal	senat staff said
richmond	said fred axelgard	scal war	senior
ridicul	said intern tanker market	scandinavia	senior figur
rift	said iranian	scar	sens
rig	said iranian forc	sceptic	sent
right	said japan	schedul	sentim
ris	said jim mcgroarti	schlesing	sept
risen	said john dosher	schuler said	septemb
risk	said kote	scof	seri
rit aid corp	said lanston jon	scott	serious
rival	said market impati	sea isl	serv
river	said pac john dosher	sea isl citi	servic
riyal	said retir adm	season	session
road deliveri	said robert brusca	season began bas	set
robert	said soviet cargo ship bound	seat	set asid plan
robert giordano	said tanker requir decreas	sec	set custom repurchas
robertson	said waseda univers professor	sec fil	set stock split
rol	mitsuru	second	settl
roll	said washington	second half	settle
rom	sail	second quarter	sever drought
ronald	sal	secret	sever flood
room	sama	secretari	shandong
root caus	sama lower	secretari immedi	shanxi
ros	sama offer	secretari jam baker	shar
rostan offshor oilfield	samjen	section	shar adjust

shar amount	shr profit ct loss	sorghum	spin
shar cash	shr profit ct loss ct	sort	spirit
shar clos	shr profit ct profit	sought	spit
shar issu	shut	sourc	split
shar tender	sichuan	sourc said	spok
sharehold	sickli chemic	south africa	spokesman
sharehold group	sid	south african maiz	spokesman said
sharp	sign	south australian	spokeswoman
sharp drop	signal	south korea	spokeswoman phylli oaklei
sharpli	signific	south west queensland	spot
sharpli higher	signific oil reserv	southern philippin	spot market
shed unrel busi	significantli	southward	spot pric
shell	silkworm	soviet	spotlight
shelter	silkworm missil strik	soviet tanker inti	spread
shelter set	similar	soviet tanker set	spur
shift	similarli	soviet union	srd
ship	simpli	soybean	stabil
ship forc	simultan	soybean cak	stabil currenc
ship sourc said	sit astrid	soybean oil	stabil currenc valu
shipment	situat	soybean total	stabl
shop	siz	spain	staff
shop center	slap mln dlr	spanish barlei	stai
shop closur	slid	spanish market	stai special long
short	slight	spark	stak
short dollar	slight review	speak	stamford
short term	slow	special	stand
short wai outsid	smack	special tender	stand athwart
shortag	small	specialti	standard
shortli	small deposit	specialti busi	standard oil
shown	small group	specialti product	stapl petrochem commod
shr	smaller volum award	specif warn	stark possibl
shr ct	smith barnei	specter	start
shr ct ct	smooth	specul	start cover
shr dilut	soft red winter wheat	specul capit inflow	stat
shr dilut ct ct	soft wheat	specul sel	stat depart renew
shr loss	soi cak	speech	stat depart sai
shr loss ct loss	sold	speech prepar	statem
shr loss ct profit	solution	spell	statem said
shr loss ct profit ct	solv	spend	station
shr primari	somebodi head	spent	statoil
shr profit	soon	spik	stead

steel	strik	surplu	telephon
steel product	strong	surpris	televi interview
steer	strong ris	suspend	televi set
step	stronger	suva	tempo
stephen marri	stronger countermeasur	swamp	temporari
sterl invest bank group	structur	swiss	temporari reserv
stg	structur remain	swiss bank corp	tender
stick	studi	swiss franc	tender expect
stimul	stun	sydnei	tender offer
stimul action	style	syracus	tension
stimul domest demand	styren	system	tepid
stimulu	subject	system corp	term
sto	submit	taipei	term debt
stock	subsidi	taiwan	term grain suppli agreem
stock dividend	subsidiari	taiwan import	termin
stock exceed	substanti	taken	territori water
stock exchang	substanti higher pric	takeov	terror
stock market	substanti increas	takeov attempt	texa
stock market reaction todai	substitut	takeov bid	texa democrat
stock ownership	subtract	takeov offer	texa democrat senat lloyd
stock ownership plan	sudden wall street retreat	talk	texa gulf coast
stock pric	suddenli	tanaka	textil
stock pric tumbl	suffer	tank trailer	thailand
stock sal	sugar	tanker	thing
stock split	suit	tanker rat	think
stock valu	sullivan	target	think rat
stockhold	sumitomo bank	tariff	think tank
stoltenberg	summer	tarnish	third
stoltenberg said	summit	tat	thou bushel
stoltenberg told	sundai	tax	thought
stood	sunflowerse	tax credit	thousand
storm	sunflowerse cak	tax cut	threat
strait	sunflowerse nil nil	tax loss	threaten
strateg	sunflowerse oil	tax loss carryforward	threaten crop
strateg invest	supermarket	teach	threaten food shortag
strategi	superpow	technolog	threshold
stream	suppli	tehran	thrust
streamlin	supplier	tehran radio	thursdai
strength	support	telecommun	thursdai declar
strengthen	support kuwait tanker	telecommun industri	thwart
stress	surfac baker look respons	telecommun market	ti

tier structur ris feet	trad deficit	turkei	unsightli mass
tighten	trad deficit drag	turn	unsucces
tighten credit	trad deficit lead	turnov	unsucces attempt
tim	trad disput	twister	unsucces bid
tim earn	trad estim	type	unusu
tit	trad friction	uae	unveil
tobacco	trad friction diplomaci	uchida	unwelcom
today	trad gap	uchida said	upcom summit
tokyo	trad gap narrow	ullman said washington	upland cotton
tokyo dealer	trad partner	ultim destin	upland domest raw cotton
tokyo open	trad row	unabl	upturn
told	trad sanction	unchang	upward pressur
told report	trad situat	unchang rat	urg
told reuter	trad sourc	understand	usda
tomorrow	trad sourc said	underwrit	usda caution
ton	trad stand	undisclos	usda discuss
tonn	trad surplu	undisclos term	usda estim australia wheat crop
tonn fob	trader	unfavor reaction	usda estim european commun
tonn south african maiz	trader said	unhappi	crop
tonn valu	transac	union	usda gav detail breakdown
top yen today	transfer	union carbid corp	usda project export
tomado	translat	union texa	usda report corn sold
toronto	transport	union texa oil reserv drop	usda set
torranc	travel	union texa said	usdaprj
total	treasuri	unit	ussr
total cie	treasuri baker	unit arab emir	ussr bui
total explor	treasuri secretari jam baker	unit kingdom	ussr nil nil nil nil
total explor area	treasuri secretari jam baker cam	unit sal	util
total explor australia pty	treasuri secretari statem help	unit stat	util market
total hold	treasuri sourc said	univers	valu
total nil nil nil nil	trigger	unload arm	vancouv
total outstand	trigger larger corn	unmil ric	veget oil
total reserv	trim	unnam	venezuela
touch	troubl	unnam destin	venic
tpn	trust	unpaid acreag reduc requir	venic econom summit
track	try	unravel	venic jun
trad	trying	unreason	ventur
trad ban	tuesdai	unrel	vessel
trad data	tulsa	unrel busi	vibranc
trad data seen	turbul ahead	unsatisfi	vibrant economi badli
trad deal	turf	unsettl	vic presid helmut schlesing

vienna	water	western world	worri
view	waterburi	westhem	wors
vigor explor program	watersid worker feder	wheat	worsen
violat	waterwai	wheat acreag reduc requir	worth
violenc	waukesha	wheat bran	wright
virtual	wav	wheat continu	writ
virul anti american polici	weak	wheat crop	wtc
visit	weaken	wheat export	wti
vital	weaker	wheat midg	year
vital interest	wednesdai	wheat product	year ago
vital suppli	week	wheat red mit	year ago period
volcker	week ago	whit wheat	year agreem sign
volcker nightmar	week dash	whitten	year earlier
volum	week end april	wholli	year end march
vot	weekend	wholli own subsidiari	year exclud
wai	weight	wi	year fall
wai street	weinberg said	wil	year iran iraq war
wak	welcom	wing	year loss
wall street	wellington	winter	year low
wall street drop	went	wip	year meet
wall street fear	west	wireless plc	year presidenti
wall street retreat	west berlin	wisconsin	year term
waltham	west berlin stat central	withdraw	year trying
want	west german	withdraw right	yel
war	west german bundesbank	withdrew	yen
ward	west german financ minist	wood	yen strength
warlik step	gerhard	worcest	yen today
warm water	west german govern	word	yen worth
warn	west german offici mondai	work	yesterdai
warrant	west german polici	worker	yield
warship	west germani	world	york
wash	west germani withdrew objec	world bank	york corp said
washington	west texa	world economi	york law
washington announc	west texa intermedi	world expan	york said
washington privat center	west texa sour	world market	young
washington seiz	western	world stock	zealand
washita	western democraci	world trad	zurich
watch	western diplomat sai iraq start	worldwid	

ภาคผนวก ค.

ผลงานวิจัยที่เกี่ยวข้องกับวิทยานิพนธ์และได้รับการตีพิมพ์

- P.Jitpakdee and W.Kreesuradej, “Dimensionality Reduction of Feature for Text Categorization”, Proceeding of The IEASTED International Conference on Advances in Computer Science and Technology (ACST 2007), Phuket, Thailand, April, 2007.

Canadian Secretariat
THE INTERNATIONAL ASSOCIATION OF SCIENCE
AND TECHNOLOGY FOR DEVELOPMENT



88, Suite #101, 2509 Dwyer Avenue SW., Calgary, AB T2D 7L9 Canada
Telephone: (403) 249-1195 Fax: (403) 247-8861
E-mail: s.calgary@iasted.com (for conferences) and journals@iasted.com (for publications)
Web Site: <http://www.iasted.org>

11/12/2006
Pansut Jitpakdee
King Mongkut's Institute of Technology Ladkrabang
Faculty of Information Technology
Ladkrabang
Bangkok
Thailand 10520

Dear Miss Jitpakdee,

Re: 559-120 Dimensionality Reduction of Features for Text Categorization

Congratulations, your paper has been accepted for oral presentation and publication at the IASTED International Conference on Advances in Computer Science and Technology (ACST 2007), to be held from Apr 02, 2007 to Apr 04, 2007, in Phuket, Thailand. We cordially invite you to attend and present your paper at the conference. We also encourage you to register and book your flight as soon as possible, if you have not already done so.

Please complete the following by the registration deadline of Feb. 12, 2007:

1. Registration Form and Payment (mandatory)
2. Author Information Form (mandatory)
3. Hotel Reservation Form (if you need assistance in reserving a hotel room)
4. Copyright Form (mandatory)

The above materials are available for viewing on our website at <http://www.iasted.org/conferences/home-559.html>. Please let me know if you have any questions regarding registration.

Once again, congratulations on your ACST 2007 acceptance. We are very excited to be able to include your research and ideas in the conference, and we look forward to seeing you in Phuket, Thailand.

Sincerely,

Amy Christenson
Conference Manager

Please note - For Visa Purposes:

The International Association of Science and Technology for Development (IASTED), is a non-profit organization founded in Zurich, Switzerland in 1977. The purpose of IASTED is to promote economic development through science and technology. After 30 years, IASTED continues to bring top scholars and members of industry together to develop and share new ideas, facilitate cultural exchange and encourage international unity.



The International Association of Science
and Technology for Development

April 2-4, 2007
Phuket, Thailand

Conference Program

The Third IASTED International
Conference on

Advances in Computer Science and Technology

SPONSOR

The International Association of Science and
Technology for Development (IASTED)

- Technical Committee on Artificial Intelligence
- Technical Committee on Computer Graphics
- Technical Committee on Databases
- Technical Committee on Parallel & Distributed
Computing & Systems
- Technical Committee on Software Engineering

LOCATION

Novotel Phuket Resort
Kallm Beach, Patong
Phuket 83150 Thailand
Tel: +66 (0)76 342 777
Fax: +66 (0)76 342 188

DIMENSIONALITY REDUCTION OF FEATURES FOR TEXT CATEGORIZATION

Parisut Jitpakdee* and Worapoj Kreesuradej†

Faculty of Information Technology

King Mongkut's Institute of Technology Ladkrabang, Bangkok, Thailand

*E-mail: jajah_pj@hotmail.com and † E-mail: worapoj@it.kmitl.ac.th

ABSTRACT

This paper proposes a new technique for dimensionality reduction of features for text categorization. Unlike conventional method, our phrase features are generated based on word sequences of different length (Multi-grams) from phrases extracted from whole documents. Then, we utilize Odds ratio (*OR*) to perform phrase feature selection. From preliminary experiments, the proposed techniques show better performance than that of conventional methods.

KEY WORDS

Data mining, text mining, text categorization

1. Introduction

With the rapid growth of online documents available, information organization and retrieval have become a great importance. An automated text classification has been utilized as a method to organize online documents. Text categorization (also known as text classification) is, simply, the automated assignment of natural language texts to predefined categories based on their content [1]. Bag of word (BOW) is the approach most commonly used to represent documents in Text Categorization (TC) [2,3]. In BOW, each document is represented by a vector that contains an importance weighting for every words in the collection. The main drawback of the BOW representation is in destruction of semantic relations between words. Many researches attempt to preserve this information by using phrases as features, the results from the study of David Lewis [4] has shown that phrase-based representation not improve performance of the categorization because it a) contain a high number of terms, b) have an uneven distribution of feature values, c) contain many redundant features, and d) contain a lot of noise.

Several efforts have been made to circumvent the possible problems posed by using phrases and some research results showed that the addition of n-grams to the BOW representation improved performance [10]. But the main problem in using n-grams is the huge potential increase in the number of features. To reduce numbers of features,

many papers extract the top-scored features using various feature selection methods [5,6,7]

In addition, unlike conventional method, our phrase features are generated based on word sequences of different length (Multi-grams) from phrases extracted from whole documents. The experiment results show that using Multi-grams and words has better performance than that of using only word features.

This paper is organized as follows: section 2 we introduce dimensionality techniques and odds ratio, section 3 describe our method and our algorithm, section 4 present the experiments and section 5 is conclusion of this paper.

2. Dimensionality Reduction Techniques

Basically, dimensionality reduction techniques can be classified into two approaches: feature generation and feature selection [1].

2.1 Feature Generation (or Feature Induction)

New features, which are not necessarily words, are sought for representation. Usually, the new features are synthesized from the original set of features.

There are two common approaches to feature generation. The first one combines features using disjunction only. In this approach features are grouped into subsets and each such subset is then considered as a new feature. Stemming and word clustering belong to this family of methods. The second approach groups features using only conjunctions, for example, by grouping consequent or close (in proximity) words into phrases. The use of n-grams is a common method in this family.

2.2 Feature Selection (or Feature Reduction)

Feature Selection studies how to select a subset of term (e.g. words) that are used to representing documents. Its purposes include reducing dimensionality, removing irrelevant and redundant features, reducing the amount of data for categorization. Feature Selection works by ranking all the terms and then selecting some percentage. A variety of ranking criteria have been used in text

categorization. One of the simplest methods is selecting terms with medium to high document frequency (*DF*). Other efficient term selection function, widely use, are χ^2 [1], information gain [1], mutual information [1]. In this paper we use Odds Ratio (*OR*) [1] measuring discriminating power for each candidate features.

Odds ratio (*OR*) measures the odds of the word occurring in the positive class normalized by that of negative class. The basic idea is that the distribution of features on the relevant documents is different from the distribution of features on the non relevant documents. It is defined as follows:

$$OR(t_i, c_i) = \log \frac{P(t | c_i)[1 - P(t | \bar{c}_i)]}{[1 - P(t | c_i)]P(t | \bar{c}_i)} \quad (1)$$

Where $P(t | c_i)$ is the conditional probability of term t occurring in category c_i , $P(t | \bar{c}_i)$ is the conditional probability of term t occurring not in category c_i .

3. Our Method

In Our method, first if they are web documents we will remove all tag (such as <html>,
, etc.), then we use each sentence in each document as initially phrase. After that we separate each sentence into phrase by remove stop words and stem every word in phrase with Porter's algorithm [8]. For each distinct phrase in each category, we generate multi-grams (set of word sequences that have varied length from 1 to size of phrase) by using sequence of words in phrases. Our multi-gram features difference from n-grams that they have various sizes and we have not to fix maximum size of gram. Table 1 shows some phrases and its multi-grams obtained from our method.

Table 1
Example Phrase and its Multi-grams

Phrase:	"efficient text categorization algorithm"
4-gram	"efficient text categorization algorithm" (<i>OR</i> = 0.5)
3-grams	"efficient text categorization" (0.4), "text categorization algorithm" (0.8)
2-grams	"efficient text" (0.2), "text categorization" (0.8), "categorization algorithm" (0.6)
1-gram	"efficient" (0.4), "text" (0.7), "categorization" (0.8), "algorithm" (0.65)

Then, we calculate Odds Ratio (*OR*) for each gram and we select only gram in each category that has *OR* more than all its sub-grams as features of their category. Finally, we union all features from every category to be as our Features set.

From example in table 1 with our method, the finally result features set will be {"efficient", "text", "categorization", "algorithm", "text categorization", "text categorization algorithm"}. The algorithm of our method is present in figure 1.

```

Input:  $D$  – training set of documents
          $C$  – set of categories
          $min\_f$  – minimum frequency
Output: Set of Features
Procedure Features Generation

Features =  $\phi$ 
For each category  $c \in C$ 
   $P = \phi$ 
   $W = \phi$ 
  For each document  $d \in D$  and  $d$  is in category  $c$ 
     $Wdoc = \phi$ 
     $PH =$  set of extracted phrases from  $d$ 
     $P = P \cup PH$ 
  End
   $MG = \phi$ 
  For each phrase  $p$  in  $P$ 
    For  $i=1$  to  $size(p)$ 
       $MG = MG \cup$  grams size  $i$  from  $p$ 
    End
  End
   $CF = \phi$ 
  For each gram  $mg$  in  $MG$  that has frequency  $\geq min\_f$ 
    Let  $sg =$  sub-grams of  $mg$ ; 1-grams to  $size(mg)-1$  grams
    If  $OR(mg, c) > \max(OR(sg, c))$ 
       $CF = CF \cup mg$ 
    End
  End
  Feature = Feature  $\cup CF$ 
End

```

Figure 1. Feature Generation Algorithm

4. Experiments

4.1 Test Collections

We use web data set, called the J-series (available at http://ftp.cs.umn.edu/dept/users/boley/PDDP_data/). The J-series contains 185 documents and ten categories showed in Table 2.

We generate multi-gram from all documents. Table 3 present number of bigrams (from all words) and number of multi-grams (from phrases) from J-series data set but not include 1-gram because we want to compare only phrase features. It shows that our method can reduces number of feature more than 70% of conventional method.

Table 2
J-series Data Set

CATEGORY	LABEL	DOCUMENTS
Affirmative action	AA	20
Business Capital	BC	19
Information Systems	CIS	19
Electronic Commerce	EC	19
Intellectual Property	IP	19
Employee Rights	LER	16
Materials Processing	MP	17
Personal Management	PM	19
Manufacturing Systems	SI	18
Industrial Partnership	ZIPT	19

Table 3
Number of Bigrams from J-series Data Set

Class	Bigram	Multi-gram (not include 1 grams)	Percent of Reduction
AA	132	38	78.79%
BC	660	156	76.36%
CIS	294	86	70.75%
EC	183	82	55.19%
IP	752	137	81.78%
LER	1,193	217	81.81%
MP	1,996	498	75.05%
PM	99	40	59.60%
SI	2,414	635	73.70%
ZIPT	467	190	59.31%
Total (distinct)	7327	2079	71.63%

4.2 Categorization

For categorization, each document is represented as Feature vectors according to the standard vector space model with *tfidf* (term frequency-inverse document frequency) weighting [9], *tfidf* is defined as:

$$tfidf(t_k, d_j) = \#(t_k, d_j) \cdot \log \frac{|Tr|}{\#Tr(t_k)} \quad (2)$$

where is $\#(t_k, d_j)$ denotes the number of time term t_k occurs in document d_j , $\#Tr(t_k)$ is the number of documents in which t_k occurs, and $|Tr|$ is the number of documents in training set. The weights resulting from *tfidf* are often normalized by cosine normalization, given by:

$$w_{kj} = \frac{tfidf(t_k, d_j)}{\sqrt{\sum_{s=1}^{|T|} (tfidf(t_s, d_j))^2}} \quad (3)$$

where w_{kj} is the weight of term k in document j , and $|T|$ is the number of all terms.

In the vector-space model, the similarity between two documents d_i and d_j is commonly measured using the cosine function [9], given by

$$S(d_i, d_j) = \frac{\sum_{k=1}^{|T|} w_{ki} \cdot w_{kj}}{\sqrt{\sum_{k=1}^{|T|} w_{ki}^2} \cdot \sqrt{\sum_{k=1}^{|T|} w_{kj}^2}} \quad (4)$$

In this paper, we use k -nearest neighbor (kNN) method for categorize documents [9]. Given a test document d , the kNN classifier finds the k nearest document among the training documents, using document-document similarity. The similarity score of each neighbor document to the test document is used as the weight of the categories pre-assigned to the training document. The second step is to estimate the likelihood of each category by summing the weight of the category of the k nearest documents as follows :

$$P(C_j|D_x) \approx \sum S(D_x, D_j) P(C_j|D_x) \quad (5)$$

where $P(C_j|D_x) \in \{0,1\}$, if $P(C_j|D_x) = 0$ mean that D_x is not respect to category C_j , and $P(C_j|D_x) = 1$ if D_x is respect to category C_j .

In our experiments, we use k values = 1, 2, 3..., 20 for kNN categorization

4.3 Performance Measure

We applied 4-cross-validation and use average accuracy rate to measure the correctness of categorization where accuracy rate is the ratio of documents classified correctly to the category c_i and other than c_i among all the documents. The result of the experiments show in table 4.

Table 4
Result of Categorization

Feature Type	Accuracy (%)
Words	83.23
Phrases	71.89
Words + Phrases	74.61
Words + Bigrams	84.30
Words + Selected Bigrams	85.38
Multi-grams (Our Method)	86.42

5. Conclusions

This paper proposes a new technique for dimensionality reduction of features for text categorization. From preliminary experiments, the proposed techniques show better performance than that of conventional methods. In the future, extensive theoretical studies and experiments will be presented.

References

- [1] F. Sebastiani, Machine Learning in Automated Text Categorization. *ACM Computing Surveys*, 34(1), March 2002, 1-47
- [2] T. Dumis, J. Platt, D. Heckerman, and M. Sahami, Inductive learning algorithms and representation for text categorization. *CIKM '98, ACM Info. & KM*. Bethesda, US, 1998, 148-155.
- [3] M. Weiss, C. Apte, F. J. Damerau, et al., Maximizaing testing performance. *IEEE Intell. Sys.* 14(2), 1999, 63-69.
- [4] D. D. Lewis, An evaluation of phrasal and clustered representations on a text categorization task. *SIGIR '92 15th*, Kobebhavn, DK, 1992, 37-50.
- [5] D. Mladeni and M. Grobelink, Word sequences as features in text learning. *17th Electro. & Comp. Science Conf.*, Slovenia, 1998, 145-148.
- [6] C. M. Tan, Y.F. Wang and C. D. Lee, The use of bigrams to enhance text categorization. *Info. Proc. & Management*, 38(4), 2002, 529-546.
- [7] R. Bekkerman and J. Allan. Using Bigrams in Text Categorization. *CIIR Technical Report IR-408*, University of Massachusetts, 2003.

- [8] M. F. Porter, An algorithm for suffix stripping. *Program*, 14(3), 1980, 130-137.
- [9] Y. Yang, An Evaluation of statistical approach to text categorization. *Information Retrieval*, 1(1/2), 1999, 69-90.
- [10] M. Yetisgen and W.Pratt, The effect of feature representation on MEDLINE Document Classification. *AMIA '05*, Washington D.C., October, 2005.

ประวัติผู้เขียน

ชื่อ-นามสกุล	นางสาวปรีสุทธิ์ จิตต์ภักดี
วัน เดือน ปีเกิด	2 พฤษภาคม 2525 ที่โรงพยาบาลแม่และเด็ก จังหวัดเชียงใหม่
ที่อยู่	41 ถ.ป่าพร้าว ต.ป่าแดด อ.เมือง จ.เชียงใหม่ 50100
ประวัติการศึกษา	ปีพุทธศักราช 2547 สำเร็จการศึกษา วิทยาศาสตร์บัณฑิต (วิทยาการคอมพิวเตอร์) จากมหาวิทยาลัยเชียงใหม่
ผลงานวิจัย	- Jeerayut Chaijaruwanich, Chalothorn Liamwirat, Parisut Jitpakdee, et al. "DNA Microarray Data Analysis System," The 15 th Annual General Meeting of the Thai Society for Biotechnology (TSB), Chiang Mai, February, 2004. - P.Jitpakdee and W.Kreesuradej, "Dimensionality Reduction of Feature for Text Categorization", Proceeding of The IEASTED International Conference on Advances in Computer Science and Technology (ACST 2007), Phuket, Thailand, April, 2007.