

ห้องสมุดคณะเทคโนโลยีสารสนเทศ พระจอมเกล้าลาดกระบัง

การหาค่าความสัมพันธ์ระหว่างคำศัพท์กับหมวดหมู่

RELATION TEST BETWEEN
ONLINE DICTIONARY VOCABULARY AND CATEGORY



H004450



อาจารย์ที่ปรึกษา

ผศ.ดร. พรฤดี เนติโสภากุล

อพ.
๑๙/๑๑
๒๕๔๙
เลขหมู่.....
เลขทะเบียน..... 04450
วัน,เดือน,ปี. - 5 ส.ย. 2551

b. 119.228.25
i.....

รายงานนี้เป็นส่วนหนึ่งของวิชาโครงการพัฒนาระบบงาน
หลักสูตรวิทยาศาสตรมหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ
คณะเทคโนโลยีสารสนเทศ

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ภาคเรียนที่ 2 ปีการศึกษา 2549
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมีเหตุดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

**RELATION TEST BETWEEN
ONLINE DICTIONARY VOCABULARY AND CATEGORY**



**A SYSTEM DEVELOPMENT PROJECT
OF THE REQUIREMENT FOR THE DEGREE OF
MASTER OF SCIENCE PROGRAM IN INFORMATION TECHNOLOGY
FACULTY OF INFORMATION TECNOLOGY
KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG**

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อ **2/2006** เท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



COPYRIGHT 2007

FACULTY OF INFORMATION TECHNOLOGY

เอกสาร **KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG** วิชาการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

หัวข้อ	การหาค่าความสัมพันธ์ระหว่างคำศัพท์กับหมวดหมู่
นักศึกษา	นางสาวเอมอร สิริศุภางค์
รหัสนักศึกษา	48066539
ปริญญา	วิทยาศาสตรมหาบัณฑิต
สาขาวิชา	เทคโนโลยีสารสนเทศ
แขนงวิชา	วิทยาการสารสนเทศ
ปีการศึกษา	2549
อาจารย์ที่ปรึกษา	ผศ.ดร. พรฤดี เนติโสภาคกุล

บทคัดย่อ

ในปัจจุบัน มี Dictionary ต่างๆมากมาย อยู่หลายรูปแบบ ทั้งแบบหนังสือ และแบบ dictionary online ซึ่งมีคำศัพท์ใหม่ๆ เกิดขึ้นอยู่เสมอ ดังนั้น การค้นหาคำศัพท์จึงอาจจะไม่พบความหมายของคำศัพท์ที่ต้องการ และขณะนี้ก็มี website ต่างๆเกิดขึ้นมากมาย ซึ่งบนหน้า website มักมีคำศัพท์ใหม่ๆเสมอ และบางครั้งก็อาจพบคำศัพท์ที่เป็นความหมายของกันและกัน และสามารถนำมาใช้ประโยชน์ได้ ในรายงานฉบับนี้จึงได้วิเคราะห์วิธีการจัดทำ ดิกชันนารีอังกฤษ-ไทยจากเว็บไซต์ ด้วยทฤษฎี Automatic Collection of Vocabulary and Related Term from the Web โดยค้นหาคำศัพท์จากหน้า website และนำมาเป็นเก็บรวบรวมไว้ และหาความสัมพันธ์ของหมวดหมู่ที่จัดเตรียมไว้ และคำศัพท์ที่เก็บรวบรวมได้ เพื่อใช้เป็นประโยชน์ในการทำ dictionary online ซึ่งจะได้คำศัพท์ใหม่ๆมากมายที่รวบรวมได้จาก website เหล่านี้ และใช้ประโยชน์จากการทราบถึงความสัมพันธ์ของคำศัพท์ต่างๆ เพื่อเป็นประโยชน์ในการสืบค้นหาความหมาย และเพื่อพัฒนาระบบงานอื่นๆ เช่น machine translation ต่อไปในอนาคต

Title	Relation test between online dictionary vocabulary and category
Student	Miss Aimon Sirisuphang
Student ID.	48066539
Degree	Master of Science
Programme	Information Science
Academic Year	2006
Advisor	Asst.Prof. Dr.Ponrudee Netisopakul

ABSTRACT

Nowadays, there are many dictionaries in many different format, e.g. book, online website. New words appear everyday, so finding the desire word's meaning maybe difficult. However, the online websites always contain new vocabularies, and they are also the meaning of the other words. These websites are very useful. This project analyzes the method of creating Dictionary English - Thai from website with Automatic Collection of Vocabulary and Related Term from the Web. This method searches vocabularies from websites and stores on the database, finds the relation between prepared categories and the stored vocabularies. The result is dictionary online which contains many new vocabularies from websites. Moreover, the relation values of each vocabulary can be used for other searching its meaning, and for other further research, such as, machine translation.

กิตติกรรมประกาศ

การจัดทำโครงการพัฒนาระบบงานในหัวข้อเรื่องการหาค่าความสัมพันธ์ระหว่างคำศัพท์ กับหมวดหมู่นี้สำเร็จได้นั้นด้วยความกรุณาจากอาจารย์ที่ปรึกษา ผศ.ดร. พรฤติ เนติโสภากุล ที่กรุณา ให้คำปรึกษาแนะนำ และช่วยตรวจสอบแก้ไขข้อบกพร่องของโครงการนี้ ตลอดจนให้ความรู้และ ข้อคิดเห็นที่เป็นประโยชน์อย่างยิ่งต่อโครงการ

ขอขอบพระคุณ คุณชัยยศ รักจิตเวชสกุล และ คุณชัชพงศ์ สุธีสุขสถาพร Software engineer บริษัท Microsoft (Thailand) ที่ช่วยสนับสนุนการทำโครงการ ช่วยให้คำแนะนำแนวทาง และข้อเสนอแนะที่เป็นประโยชน์ยิ่งต่อทำโครงการ

และขอขอบคุณเพื่อนๆ ทุกคนที่ให้กำลังใจซึ่งกันและกัน แลกเปลี่ยนประสบการณ์ซึ่งกัน และกัน จนทำให้การพัฒนาโครงการนี้สำเร็จลุล่วงไปด้วยดี

สุดท้ายนี้ข้าพเจ้าขอกราบขอบพระคุณ บิดา มารดา และครอบครัวของข้าพเจ้าที่ร่วมเป็น กำลังใจ และให้การสนับสนุนในทุกๆ เรื่อง ด้วยดีเสมอมา

สำหรับคุณงานความดีและประโยชน์อันพึงมาจาก โครงการฉบับนี้ ข้าพเจ้าขอมอบให้กับ บิดามารดา ซึ่งเป็นที่รักและเคารพยิ่ง ตลอดจนครูอาจารย์ที่เคารพทุกท่านที่ได้ประสิทธิ์ประสาทวิชา ความรู้และถ่านทอดประสบการณ์ที่ดีให้แก่ข้าพเจ้า

เอมอร สิริศุภางค์

สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	I
บทคัดย่อภาษาอังกฤษ.....	II
กิตติกรรมประกาศ.....	III
สารบัญ.....	IV
สารบัญตาราง.....	VI
สารบัญรูป.....	VII

บทที่ 1. บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา.....	1
1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา.....	2
1.3 สมมติฐานของการศึกษา.....	2
1.4 ทฤษฎีหรือแนวคิดที่ใช้ในการวิจัย.....	3
1.5 การเปรียบเทียบระหว่างวิธีการที่นำเสนอกับวิธีการแบบพื้นฐาน.....	3
1.6 ขอบเขตการวิจัย.....	4
1.7 ขั้นตอนของการศึกษา.....	4

บทที่ 2. ทฤษฎีที่เกี่ยวข้อง

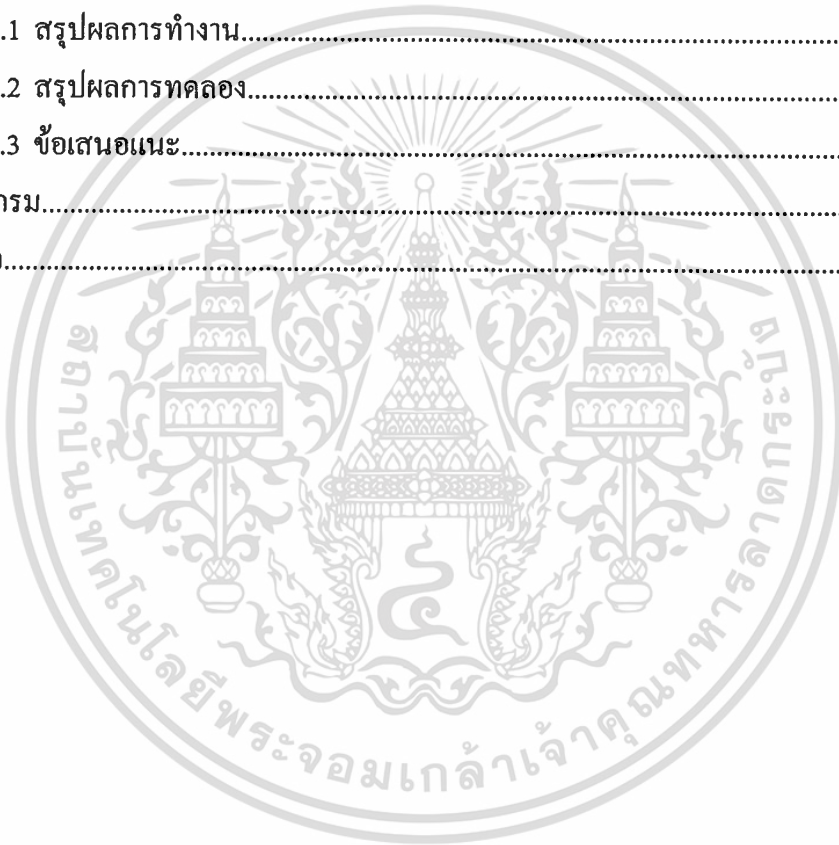
2.1 ทฤษฎีการทำ Automatic Collection of Vocabulary and Related Term from the Web.....	6
2.2 กระบวนการทำ Automatic Collection of Vocabulary and Related Term from the Web.....	6
2.3 งานหลักของการทำ Automatic Collection of Vocabulary and Related Term from the Web.....	8

บทที่ 3. วิธีการดำเนินการศึกษา

3.1 โครงสร้างการทำงานตามทฤษฎี.....	11
3.2 ขั้นตอนการดำเนินงาน.....	11
3.3 การพัฒนาโครงการตามลำดับขั้น.....	25
3.4 รูปแบบของผลลัพธ์ที่ตามขั้นตอนการทำงาน.....	29

สารบัญ (ต่อ)

	หน้า
บทที่4. ผลการทดลอง	
4.1 การเตรียมข้อมูล.....	30
4.2 ส่วนของการทำงานของ โปรแกรม.....	32
บทที่5. สรุปผลและข้อเสนอแนะ	
5.1 สรุปผลการทำงาน.....	45
5.2 สรุปผลการทดลอง.....	46
5.3 ข้อเสนอแนะ.....	46
บรรณานุกรม.....	47
ภาคผนวก.....	48



สารบัญตาราง

ตารางที่	หน้า
4.1 แสดงข้อมูลที่ใช้ในการทดลองของ Directory.....	31
4.2 แสดงข้อมูลที่ใช้ในการทดลองของ Subdirectory.....	31
4.3 แสดงข้อมูลที่ใช้ในการทดลองของ Dictionary.....	32
4.4 แสดงผลที่ได้จากการทดลองของ Relation test.....	37



สารบัญรูป

รูปที่	หน้า
2.1 แสดงการทำ Automatic Collection of Vocabulary and Related Term from the Web โดยสรุป.....	7
3.1 เปอร์เซนต์ของเวลาที่ใช้แต่ละขั้นตอน	13
3.2 ลักษณะของข้อมูลของแต่ละ File ใน corpus ที่รวบรวมได้จากการทำ Web page collection.....	14
3.3 ลักษณะของข้อมูลของแต่ละ File ใน corpus ที่รวบรวมได้จากการทำ Sentence extraction.....	15
3.4 ลักษณะของข้อมูลของแต่ละ File ใน corpus ที่รวบรวมได้จากการทำ Automatic term recognition.....	15
3.5 แสดงรูปแบบของฐานข้อมูล (Database diagram).....	16
3.6 แสดงวิธีการเก็บหมวดหมู่และหมวดหมู่ย่อยลงในฐานข้อมูล.....	19
3.7 แสดงวิธีการนับจำนวน word hits จาก search engine.....	20
3.8 แสดงวิธีการนับจำนวน both hits จาก search engine.....	21
3.9 แสดงการเก็บค่าต่างๆที่ได้จากการทำ Relation test.....	23
3.10 แสดงขั้นตอนการดำเนินงาน.....	24
4.1 หน้าจอหลักของ Application.....	33
4.2 ผลลัพธ์จากการกดปุ่ม Compiling corpus ที่ปรากฏบน application.....	34
4.3 ผลลัพธ์จากการกดปุ่ม Filtering ที่ปรากฏบน application.....	35
4.4 ผลลัพธ์จากการกดปุ่ม Relation ที่ปรากฏบน application.....	36
4.5 แสดงผลที่ได้จากการทำ Relation test แบบเจาะจงคำศัพท์.....	38
4.6 แสดงรูปแบบหน้าการทำงานหลักของ online dictionary.....	39

สารบัญรูป(ต่อ)

รูปที่	หน้า
4.7 แสดงการทำงานของ การตรวจสอบ Relation test ตามหมวดหมู่.....	40
4.8 แสดงการทำงานของ การตรวจสอบ Relation test ตามหมวดหมู่ (เพิ่มเติม).....	40
4.9 แสดงการทำงานของ การตรวจสอบ Relation test ตามคำศัพท์.....	41
4.10 แสดงการทำงานของ การตรวจสอบ Relation test ตามคำศัพท์ (เพิ่มเติม).....	41
4.11 แสดงผลของ Relation test สำหรับคำศัพท์ที่มีความเกี่ยวข้องกับหมวดหมู่หมวดหมู่มาก	42
4.12 แสดงผลของ Relation test สำหรับ คำศัพท์ที่มีความเกี่ยวข้องกับหมวดหมู่.....	42
4.13 แสดงผลของ Relation test สำหรับคำศัพท์และหมวดหมู่อาจเกี่ยวข้องหรือไม่เกี่ยวข้อง กัน.....	42
4.14 แสดงผลของ Relation test สำหรับคำศัพท์และหมวดหมู่ไม่มีความเกี่ยวข้องกันเลย.....	43



บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

ในปัจจุบันความก้าวหน้าทางเทคโนโลยีด้านสารสนเทศได้รับการพัฒนาขึ้นอย่างมาก ทำให้ข้อมูลต่างๆจำเป็นต้องถูกรวบรวมไว้ในหน่วยความจำของคอมพิวเตอร์ และนำมาใช้งานได้ตรงความต้องการของผู้ใช้งานที่สุด หนึ่งในข้อมูลที่สำคัญคือคำศัพท์และคำแปล ซึ่งเป็นสิ่งที่เราสามารถนำมาใช้ประโยชน์ต่างๆได้ โดยเฉพาะการทำ Online Dictionary ในขณะนี้มี Dictionary ต่างๆมากมาย อยู่หลายรูปแบบ ทั้งแบบหนังสือ และแบบ Online Dictionary ซึ่งมีคำศัพท์ใหม่ๆ เกิดขึ้นอยู่เสมอ ดังนั้น การค้นหาคำศัพท์จึงอาจจะไม่พบความหมายของศัพท์ที่ต้องการ Dictionary คือ สิ่งที่รวบรวมการแปลภาษาจาก ภาษาหนึ่ง ไปยังอีกภาษาหนึ่ง เพื่อช่วยผู้ที่ต้องการแปลภาษาที่ตนเอง ไม่นัด เพื่อใช้ประโยชน์ต่างๆกัน ในปัจจุบันมีการทำอยู่ในหลายรูปแบบเช่น การทำเป็นหนังสือ, เครื่อง talking dictionary หรือ ที่นิยมกันมาในปัจจุบัน คือ Online Dictionary ซึ่งมีการจัดทำผ่านระบบ website เพราะผู้ใช้จะสะดวกต่อการใช้งาน เพียงติดต่อบริษัท internet เท่านั้น และขณะนี้ก็มี website ต่างๆเกิดขึ้นมากมาย ซึ่งบนหน้า website มักมีคำศัพท์ใหม่ๆเสมอ และบางครั้งก็อาจพบคำศัพท์ที่เป็นความหมายของกันและกัน และสามารถนำมาใช้ประโยชน์ได้ ในรายงานฉบับนี้จึงได้วิเคราะห์วิธีการจัดทำ dictionary อังกฤษ-ไทย โดยการรวบรวมคำศัพท์จาก website ต่างๆที่มีคำศัพท์และคำแปล นำมาเป็นเก็บรวบรวมไว้ โดยที่มีการจัดทำหมวดหมู่ต่างๆของคำศัพท์ เพื่อหาความเกี่ยวข้องกันของคำศัพท์ต่างๆ และบ่งบอกว่าเป็นคำศัพท์ที่เกี่ยวข้องกับเรื่องใด โดยที่มีการจัดเรียงคำศัพท์ที่ได้ตามหมวดหมู่ และเก็บไว้เรียงตามหมวดหมู่โดยที่จะมีค่าความเป็นไปได้ในการที่คำศัพท์แต่ละตัวจะอยู่หมวดหมู่ต่างๆอีกด้วย เพื่อใช้ให้เป็นประโยชน์ต่างๆ โดยเฉพาะการทำ Online Dictionary ซึ่งจะได้คำศัพท์ใหม่ๆมากมายที่รวบรวมได้จาก website เหล่านี้ เพื่อเป็นประโยชน์ในการสืบค้นหาความหมายต่อไปในอนาคต

การพัฒนาโปรแกรมรวบรวมคำศัพท์จากหน้าเว็บไซต์ มีการทำงานโดยใช้ Automatic Term Recognition เพื่อทดสอบความเกี่ยวเนื่องกันของคำศัพท์และหมวดหมู่การทำงาน และมีการวัดผลค่าที่ได้ว่าเป็นความเกี่ยวเนื่องกันแค่ไหน โดยใช้การทำ Relation test และเก็บผลของการทำงาน โดยใช้หลักการวิเคราะห์ความแม่นยำของผลที่ได้จากการทำ Relation test ว่าต้องปรับปรุงแก้ไขอย่างไร เมื่อนำมาใช้จริง และการวิเคราะห์ผลลัพธ์ รวมถึงการเพิ่มประสิทธิภาพในการค้นหาศัพท์ และความถูกต้องแม่นยำของศัพท์ที่ได้ เพื่อนำมาใช้ประโยชน์ในการทำ Online Dictionary และการใช้ประโยชน์เพิ่มเติมด้านอื่นๆในอนาคต

1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา

การศึกษาโครงการนี้มีวัตถุประสงค์คือ

1. เพื่อศึกษาการทำการเก็บรวบรวมคำศัพท์เป็นจากหน้าเว็บไซต์ ศึกษาวิธีการหาคำศัพท์จากเว็บไซต์พร้อมทั้งคำแปลและเก็บรวบรวม
2. เพื่อพัฒนาโปรแกรมรวบรวมคำศัพท์จากหน้าเว็บไซต์ด้วยวิธีการ Compiling Corpus , Automatic Term Recognition และหาความสัมพันธ์ของคำศัพท์และหมวดหมู่ต่างๆ Relation test (Filtering)
3. เพื่อศึกษาความถูกต้องและความแม่นยำของวิธีการทำ Relation test (Filtering) ว่าค่าความน่าจะเป็นที่ได้มีความถูกต้องมากเพียงใดเพื่อสามารถนำไปใช้ในการหาความสัมพันธ์ของข้อมูลอื่นๆ ได้ หรือควรปรับปรุงต่อไปในอนาคต
4. เพื่อประยุกต์คำศัพท์ที่ได้มาใช้ในการทำประโยชน์ด้านต่างๆ โดยที่โครงการนี้จะประยุกต์ใช้ในการทำ Online Dictionary
5. เพื่อให้คำศัพท์ที่รวบรวมได้นำไปทำโครงการอื่นได้ในอนาคตเช่น การทำ machine translation

1.3 สมมติฐานของการศึกษา

ข้อดีของการหาค้นหาคำศัพท์จากเว็บไซต์และความสัมพันธ์ของคำศัพท์และหมวดหมู่คือ ค่าความน่าจะเป็นที่คำศัพท์จะเกี่ยวข้องกับหมวดหมู่ที่เราหาได้จากสูตร อาจจะไม่ได้ออกคือ 100% ซึ่งด้วยสมมติฐานของผลลัพธ์จากวิธีการหา related term นี้ คำศัพท์ที่ได้มากกว่า 80% จะมีความเกี่ยวข้องกับหมวดหมู่นั้น แต่อาจจะเป็นส่วนหนึ่งของหมวดหมู่นั้นๆหรือไม่ก็ได้ การดึงคำศัพท์จากหน้าเว็บไซต์ และการหาค่าความสัมพันธ์ใช้เวลานานต่อศัพท์ 1 คำ และทั้งนั้นเรามีศัพท์ทั้งหมดที่จะรวบรวมและนำมาหาความสัมพันธ์มากกว่า 100,000 ตัว โดยที่มีหมวดหมู่หลักและหมวดหมู่ย่อยรวมกันมากกว่า 50 หมวดหมู่ และเราต้องหาความสัมพันธ์ของคำศัพท์ทั้งหมดกับหมวดหมู่ที่เราต้องการทราบ จะเห็นได้ว่าเป็นปริมาณจำนวนมากที่คอมพิวเตอร์จะต้องนำมาประมวลผล รวมถึงการที่ต้องมีการทำงานกับฐานข้อมูลทั้งการอ่านและการเขียนจำนวนมากกว่า 100000 ครั้ง ดังนั้นทำให้สมรรถนะของเครื่องคอมพิวเตอร์ต่ำลงซึ่งเป็นผลทำให้การทำงานของโปรแกรมช้าลงไปตามระยะเวลา ทำให้เวลานานมากในการประมวลผลทั้งการหาคำศัพท์และการหาความสัมพันธ์ ดังนั้นจึงต้องหาแนวทางในการแก้ปัญหาในจุดนี้ เมื่อลองแล้วว่าการทำงานมีปัญหาตรงจุดนี้จริง เราจึงต้องทำการดึงคำศัพท์จากหน้าเว็บไซต์จะต้องทำการตั้งเวลาให้หยุดการทำงานและเริ่มใหม่ที่จุดเดิมเป็นช่วงๆเพื่อให้การทำงานมีประสิทธิภาพดีขึ้น

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1.4 ทฤษฎีหรือแนวคิดที่ใช้ในการวิจัย

เราใช้เทคนิคการทำ Automatic Collection of Vocabulary and Related Term from the Web ที่ประกอบไปด้วยการทำ 3 เทคนิค คือ Compiling Corpus, Automatic Term Recognition, Relation test (Filtering) จากแนวคิดที่ต้องการจัดเรียงคำศัพท์เป็นหมวดหมู่แทนการจัดเรียงคำศัพท์ตามตัวอักษร และความสัมพันธ์ของคำศัพท์กับคำศัพท์และหมวดหมู่ต่างๆ โดยแนวคิดนี้จะทำให้เราสามารถดึงประโยชน์จากข้อมูลบนเว็บไซต์และ search engine website ให้เป็นประโยชน์ และเป็นแนวคิดรูปแบบใหม่ในการรวบรวมและจัดเก็บคำศัพท์ เพื่อให้สามารถนำคำศัพท์ที่เราได้มาใช้ให้เกิดประโยชน์ในแนวทางใหม่ๆ และรวบรวมคำศัพท์ที่ได้มาพร้อมค่าความสัมพันธ์มาจัดเป็นฐานข้อมูล ค่อยจะมีการจัดทำ online dictionary ที่มีรูปแบบใหม่ ซึ่งสามารถจะค้นหาความหมายของคำศัพท์ได้ในแบบเดิมๆ และในแบบใหม่ คือคำศัพท์ทั้งหมดที่มีความเกี่ยวข้องกัน หรือเกี่ยวข้องกับหมวดหมู่ โดยที่จะมีการจัดทำฐานข้อมูลใหม่ที่มีการจัดเรียงตามหมวดหมู่ เพื่อสามารถนำไปประยุกต์ใช้ในโครงการอื่นๆ ในอนาคตต่อไปได้อีกด้วย เช่นการทำ machine translator เป็นต้น

1.5 การเปรียบเทียบระหว่างวิธีการที่นำเสนอกับวิธีการแบบพื้นฐาน

จากวิธีการเก็บรวบรวมคำศัพท์แบบพื้นฐานที่เรียงตามอักษร ซึ่งการค้นหาคำศัพท์มาใช้งาน เช่นการทำ online dictionary ก็จะได้เฉพาะความหมายของคำนั้นๆ แต่ในกรณีนี้วิธีการทำ Automatic Collection of Vocabulary and Related Term from the Web ที่เสนอในรายงานฉบับนี้ จะทำให้ได้ฐานข้อมูลที่มีการจัดลำดับความสัมพันธ์ของคำศัพท์ไว้กับหมวดหมู่ต่างๆ ว่าศัพท์คำใดมีความเกี่ยวข้องกับหมวดหมู่ใดมากน้อยเพียงใด ทำให้ได้ค่าของความน่าจะเป็นในการเกี่ยวข้องกันของคำศัพท์และคำศัพท์ หรือของคำศัพท์และหมวดหมู่ ทำให้เมื่อมีการค้นหาคำศัพท์ผ่านทาง online dictionary สามารถค้นหาคำศัพท์ในหมวดหมู่เดียวกันกับคำศัพท์ที่เราต้องการค้นหาด้วยการทำ dictionary โดยใช้คำศัพท์จาก website ก่อให้เกิดประโยชน์ในการรวบรวมและค้นหาคำศัพท์ใหม่ๆ ที่เกิดขึ้นอยู่ตลอดเวลา ซึ่ง dictionary แบบเก่าที่มีมา มักไม่มีการรวบรวมคำศัพท์ใหม่ๆ หรือคำศัพท์เฉพาะด้านต่างๆเอาไว้ ทำให้การค้นหาของผู้ใช้งานเป็นไปได้โดยไม่ครบถ้วน ดังนั้น รายงานฉบับนี้จึงกล่าวถึงการจัดทำ website Online Dictionary ซึ่งเป็นรูปแบบใหม่ของการค้นหาศัพท์ โดยวิธีการที่ได้วิเคราะห์ในรายงานฉบับนี้ ทั้งหารหาคำศัพท์ การหา website การเก็บรวบรวม จนถึงการนำมาใช้งาน สามารถนำไปใช้ได้จริง ทำได้ไม่ยาก และมีประสิทธิภาพ โดยที่สามารถทำให้ระบบกลับมาค้นหาคำศัพท์อีกครั้งได้ใหม่ ตามความถี่ที่เราตั้งไว้ ซึ่งจะช่วยให้ได้คำศัพท์ใหม่ๆ เสมอ และจะเป็นประโยชน์กับผู้ใช้งาน ในการค้นหาคำศัพท์ โดยเฉพาะในการใช้คำศัพท์จาก website Online Dictionary ซึ่งจะมีคำศัพท์ใหม่ๆ และคำศัพท์แตกต่างกันไปในแต่ละ website ทำให้เราได้ข้อมูลมากมายที่รวบรวมมาจากหลายๆ website และเมื่อทำให้เกิดการ update ตามระยะเวลาที่ตั้งไว้ จะทำ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาติให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ให้เมื่อ website ที่เรารวบรวมข้อมูลมา มีการเพิ่มคำศัพท์ เราก็จะได้คำศัพท์เหล่านั้นมาด้วย ทำให้ได้คำศัพท์จำนวนมากมาย และนำมาใช้ประโยชน์ได้หลายอย่าง ซึ่งเป็นประโยชน์ต่อผู้ที่ต้องการค้นหาคำศัพท์เป็นอย่างดี

1.6 ขอบเขตการวิจัย

ผลของการศึกษาทำให้ทราบว่า มีกระบวนการการทำงานที่สามารถรวบรวมคำศัพท์และหาความสัมพันธ์ของคำศัพท์และหมวดหมู่ต่างๆ ได้ โดยรายงานชุดนี้ จะศึกษาการแปลภาษาอังกฤษ ไปเป็นภาษาไทย และภาษาไทยเป็นภาษาอังกฤษ โดยที่จะมีความแตกต่างกับ Dictionary ทั่วไป คือเป็นคำศัพท์ที่หามาจากหน้าของ website โดยการเปรียบเทียบศัพท์ อังกฤษ-ไทยที่ปรากฏบนหน้า website ซึ่งจะทำได้คำศัพท์ใหม่ๆ ที่เกิดขึ้น เนื่องจาก website มีการเกิดใหม่ทุกวัน และมักมีศัพท์ใหม่ๆ นำมาใช้งาน ดังนั้น ผู้ที่ค้นหาคำศัพท์จาก dictionary มาตรฐานที่มีอยู่เดิม แล้วไม่พบนั้น สามารถลองค้นหาจาก Online Dictionary ที่จัดทำจาก ศัพท์ในหน้า website ได้ ทำให้มีโอกาที่จะพบความหมายของคำที่ต้องการเพิ่มขึ้น

ในเอกสารสัมมนาฉบับนี้ ค้นหาวิธีการดึงข้อมูลจากหน้า Website และวิธีการเปรียบเทียบว่า คำศัพท์ใด เป็นความหมายของคำศัพท์ใด โดยใช้หลักการจากการรวบรวมข้อมูลว่ามีเงื่อนไขใดบ้างที่จะเป็นสิ่งบ่งบอกว่า คำศัพท์ใด เป็นความหมายของคำศัพท์ใด และนำคำศัพท์ที่ได้เก็บรวบรวม และหาความสัมพันธ์ของคำศัพท์และหมวดหมู่ต่างๆ โดยจัดทำเป็น Online Dictionary บน website ต่อไปในอนาคต และใช้ ทฤษฎีการทำ Automatic Collection of Vocabulary and Related Term from the Web

โดยมีขอบเขตหลักๆ ที่จะศึกษาต่อไปนี้

- การทำ Compiling Corpus
- การทำ Automatic Term Recognition
- การทำ Filtering
- การทำ Relation test

1.7 ขั้นตอนของการศึกษา

วิทยานิพนธ์ฉบับนี้ได้แบ่งเนื้อหาออกเป็น 5 บทด้วยกันคือ

บทที่ 1 กล่าวถึงความเป็นมาของงานวิจัย ความมุ่งหมายและวัตถุประสงค์ สมมติฐาน ทฤษฎีที่ใช้ ขอบเขตของการวิจัย และขั้นตอนการศึกษา

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 2 กล่าวถึงวิธีการทำ Automatic Collection of Vocabulary and Related Term from the Web และวิธีการทำ Compiling Corpus Automatic, วิธีการทำ Term Recognition, วิธีการทำ Filtering, วิธีการทำ Relation test ทฤษฎีและแนวคิดการประเมินของผลที่ได้รับจากโครงการ สำหรับแต่ละขั้นตอนของการดำเนินการ

บทที่ 3 กล่าวถึงวิธีการดำเนินการทำโครงการ เริ่มตั้งแต่การเตรียมข้อมูลเพื่อใช้ค้นหา คำศัพท์จากเว็บไซต์ การหาคำศัพท์และความหมายจากเว็บไซต์ การเก็บรวบรวมคำศัพท์ การคัดเลือกหมวดหมู่เพื่อนำมาค้นหาความสัมพันธ์กับคำศัพท์ การค้นหาความสัมพันธ์ของคำศัพท์และหมวดหมู่โดยการเก็บค่าความน่าจะเป็นของการเกี่ยวข้องกัน การนำคำศัพท์มาเรียงลงฐานข้อมูล โดยเรียงตามความสัมพันธ์ การสร้าง online dictionary เพื่อนำข้อมูลที่ได้มามาใช้ให้เกิดประโยชน์สูงสุด เพื่อใช้ในการค้นหาคำศัพท์และความหมายของคำศัพท์ต่างๆ

บทที่ 4 กล่าวถึงผลการทดลอง โดยสรุปผลของการทำการใช้ทฤษฎี Automatic Collection of Vocabulary and Related Term from the Web ในการทำโครงการ ว่าได้ผลถูกต้องแม่นยำเพียงใด ควรมีการปรับปรุงด้านใด และเหมาะนำไปใช้ประโยชน์ หรือพัฒนาโครงการใดต่อไปในอนาคต

บทที่ 5 เป็นบทสรุปผลการวิจัยและข้อเสนอแนะ

บทที่ 2

ทฤษฎีที่เกี่ยวข้อง

2.1 ทฤษฎีการทำ Automatic Collection of Vocabulary and Related Term from the Web (Sato and Y.Sasaki. 2003. “Automatic collection of related terms from the web.” 121-124. InProc.)

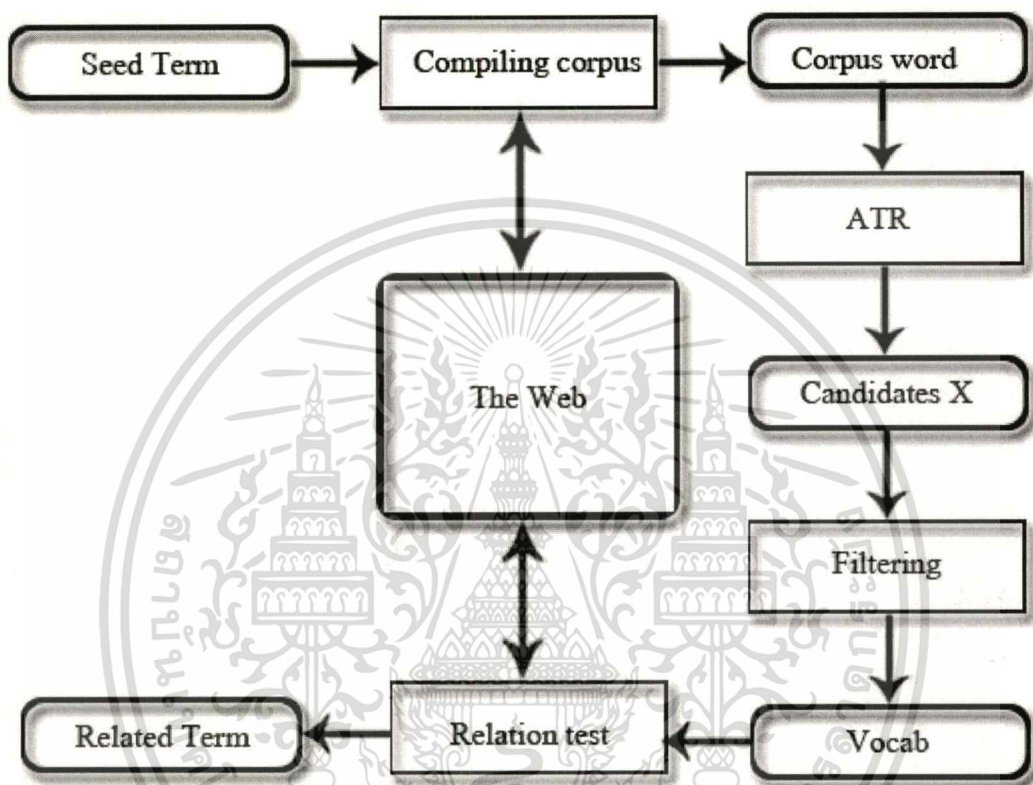
Automatic Collection of Vocabulary and Related Term from the Web คือกระบวนการในการค้นหาคำศัพท์และคำแปลจากหน้าเว็บไซต์แบบอโตเมติก และคำศัพท์ที่เกี่ยวข้องกัน โดยแบ่งตามหมวดหมู่ และหาความน่าจะเป็นในหมวดหมู่นั้นๆ แล้วรวบรวมเก็บไว้เป็น Small Corpus ทำให้ได้ข้อมูลที่เป็นคำศัพท์ใหม่ๆ จากหน้าเว็บไซต์เพื่อนำมาสนับสนุนการทำประโยชน์ในด้านต่างๆ เช่น การทำ Online Dictionary กระบวนการนี้เป็นขั้นตอนในการดึงเอาคำศัพท์และคำแปลมาจากเว็บไซต์ทำให้ได้ความคำศัพท์ใหม่ๆ ที่ไม่เคยรู้มาก่อน หรือได้คำศัพท์ที่มีความถูกต้อง หรือคำศัพท์ที่นำไปใช้ประโยชน์ได้ Automatic Collection of Vocabulary and Related Term from the Web เป็นส่วนหนึ่งของกระบวนการการค้นหาพื้นฐานความรู้จากเว็บไซต์ การนำแนวโน้มของคำศัพท์ข้อมูลต่างๆ ที่ซ่อนอยู่ในหน้าเว็บไซต์ออกมาใช้งานเป็นสิ่งสำคัญ เพราะถ้าไม่รู้จักใช้ประโยชน์จากข้อมูลเหล่านั้นก็สูญเปล่า ดังนั้น กระบวนการนี้เป็นเครื่องมือที่สำคัญในการค้นหาคำศัพท์จากเว็บไซต์ที่มีขนาดใหญ่ทำให้ได้คำศัพท์ที่มีประโยชน์ซึ่งซ่อนอยู่ในเว็บไซต์ เป็นการเพิ่มคุณค่าให้กับข้อมูลบนเว็บไซต์ที่เรามีอยู่

นอกจากนั้น Automatic Collection of Vocabulary and Related Term from the Web ยังเป็นกระบวนการที่ช่วยแบ่งคำศัพท์ออกตามหมวดหมู่ และศึกษาความน่าจะเป็นในการเกี่ยวเนื่องกันของคำศัพท์และหมวดหมู่ต่างๆ โดยเก็บข้อมูลที่ได้ไว้ในฐานข้อมูล โดยเรียงตามหมวดหมู่ของคำศัพท์ ดังนั้นกระบวนการนี้ช่วยทำให้เกิดศักยภาพในการใช้ข้อมูลในเว็บไซต์ให้มากยิ่งขึ้น

2.2 กระบวนการทำ Automatic Collection of Vocabulary and Related Term from the Web (Sato and Y.Sasaki. 2003. “Automatic collection of related terms from the web.” 121-124. InProc.)

Automatic Collection of Vocabulary and Related Term from the Web มีการรวบรวมเทคนิคต่างๆ เช่น การทำ Compiling Corpus โดยแบ่งเป็น Web page collection และ Sentence extraction, การทำ Automatic term recognition , การทำ Filtering , การทำ Relation Test และสุดท้ายคือการวิเคราะห์ Experiments และ Discussion ขั้นตอนทั้งหมดนี้ได้ถูกทำเพื่อหาคำศัพท์ที่สามารถนำไปใช้งานได้จริง และหาความสัมพันธ์ของคำศัพท์นั้นๆ เพื่อใช้ให้เกิดประโยชน์สูงสุด เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่นิยามให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

โดยทั่วไปเมื่อกล่าวถึงการทำ Automatic Collection of Vocabulary and Related Term from the Web ส่วนใหญ่จะคำนึงถึงและให้ความสำคัญกับขั้นตอนการทำ Filtering และ การทำ Relation Test เนื่องจากเป็น 2 ขั้นตอนที่ทำให้เกิดกระบวนการหาคำศัพท์ที่เกี่ยวข้องกันของ และเป็นกระบวนการที่ช่วยค้นพบข้อมูลที่เป็นความรู้จากหน้าเว็บไซต์โดยแบ่งขั้นตอนของการทำ Automatic Collection of Vocabulary and Related Term from the Web ได้ดังนี้



รูปที่ 2.1 Automatic Collection of Vocabulary and Related Term from the Web โดยสรุป

การทำ Automatic Collection of Vocabulary and Related Term from the Web เป็นการทำการค้นหาและตรวจสอบความเกี่ยวข้องกันคำศัพท์สำหรับแต่ละหมวดหมู่ของคำศัพท์ ซึ่งจะประกอบด้วยหลายขั้นตอนที่บางขั้นตอนมีการทำซ้ำๆ หรือต้องมีการวนกลับมาทำซ้ำใหม่ ซึ่งจะประกอบด้วยการทำงาน 4 กระบวนการ ได้แก่

- Compiling corpus
- Automatic Term Recognition
- Filtering
- Relation test

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.3 งานหลักของการทำ Automatic Collection of Vocabulary and Related Term from the Web แบ่งออกเป็น 4 ส่วน ดังนี้

- Compiling corpus รวบรวมหน้าเว็บไซต์ที่มีคำศัพท์และคำแปลมาจัดเก็บเป็น Corpus
- Automatic Term Recognition กำหนดเงื่อนไขในการค้นหาคำศัพท์และตัดประโยคที่มีคำศัพท์พร้อมคำแปลมาจัดเก็บไว้เป็น Corpus ที่มีขนาดเล็กลง
- Filtering ตัดคำศัพท์และคำแปลเฉพาะที่ต้องการจริงๆและนำมาเก็บลงฐานข้อมูล
- Relation test หาความสัมพันธ์ของคำศัพท์และหมวดหมู่ต่างๆและนำมาเก็บลงฐานข้อมูล โดยเรียงคำศัพท์ตามหมวดหมู่และบอกค่าความน่าจะเป็นในการอยู่ในหมวดหมู่นั้นๆ

Compiling corpus

Compiling corpus เป็นกระบวนการที่เริ่มจากการส่งค่า seed term เข้าไปใน compiling corpus โดยในที่นี้จะใช้ seed term คือ URL ของ Website dictionary และกระบวนการทำ Compiling corpus จะทำการเลือก Web page ทั้งหมดที่มีเก็บค่า seed term เอาไว้ ในที่นี้เราใช้ The Web สำหรับการดึงหน้าเว็บไซต์ที่ต้องการ และหน้าหน้าเว็บไซต์ที่เราเก็บได้จาก Compiling corpus นำมาตัดเฉพาะ passages ที่เกี่ยวข้องกับ seed term โดยการผ่านกระบวนการทั้งหมด 2 กระบวนการดังนี้

1. Web page collection คือการเก็บรวบรวมหน้า web page โดยการเอา URL ที่ได้มาทำการเปิดหน้า web page เพื่อรวบรวมคำศัพท์ และทำ compiling corpus คือเอา web page ทั้งหมดที่ได้จากการเปิด URL มาเป็น corpus ของคำศัพท์
2. Sentence extraction คือการตัดเอาแต่ละประโยคที่เราต้องการ มาจาก corpus ที่เราเก็บไว้ จะได้ออกมาเป็นประโยคที่มีคำศัพท์และคำแปลอยู่ข้างใน

Automatic Term Recognition

Automatic term recognition เป็นกระบวนการการสร้างเงื่อนไขในการตัดเอาเฉพาะคำศัพท์และคำแปลที่เราต้องการออกจาก corpus สามารถทำได้โดยการสร้าง term list ในการดึงคำศัพท์ออกจากหน้า website แบบ automatic เริ่มจากการสร้าง term list เป็นเงื่อนไขในการดึงคำศัพท์จากแต่ละหน้าเว็บไซต์และนำมาใช้ในการดึงคำศัพท์และคำแปล สิ่งที่ได้จากกระบวนการนี้จะได้เป็นชุดของคำที่มีคำศัพท์และคำแปลอยู่ข้างใน

Filtering

Filtering เป็นกระบวนการในการตัดคำที่ต้องการ ทั้งคำศัพท์และคำแปล เพื่อนำมาเก็บรวบรวมลงฐานข้อมูล โดยผลจากกระบวนการนี้คือการได้คำศัพท์มาเก็บลงฐานข้อมูล โดยการสร้างเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เงื่อนไขของการตัดคำศัพท์และนำมาใช้ในการตัดคำศัพท์และคำแปลแบบ automatic ซึ่งจะมีรูปแบบสำหรับแต่ละ Corpus แบ่งตาม web page ที่ได้รวบรวมมา เมื่อได้คำศัพท์และคำแปล ก็จะทำนำมาเก็บรวบรวมลงฐานข้อมูล เพื่อใช้ในการทำ Relation test ต่อไป

Relation test

Relation test เป็นกระบวนการในการหาความสัมพันธ์ของ หมวดหมู่ที่เราต้องการกับ คำศัพท์ในฐานข้อมูลของเรา การแบ่งหมวดหมู่ของคำศัพท์ในที่นี้เมื่อเราได้คำศัพท์และคำแปลมา เก็บไว้ในฐานข้อมูลแล้ว เราจะทำการแบ่งคำศัพท์ออกมาเป็นหมวดหมู่ และสร้างฐานข้อมูลใหม่ โดยการเรียงคำศัพท์ตามหมวดหมู่ และระบุค่าความน่าจะเป็นของการเกี่ยวข้องกันของคำศัพท์แต่ละตัวกับหมวดหมู่หลักและหมวดหมู่ย่อย

การหาความสัมพันธ์ของหมวดหมู่ที่เรากำหนดไว้และคำศัพท์ในฐานข้อมูล เราจะตรวจสอบว่าทั้ง 2 คำนี้เกี่ยวข้องกันจริงไหม โดยการเกี่ยวข้องกัน จะต้องประกอบด้วย 2 อย่าง

1. หมวดหมู่ต้องมีความหมายแคบหรือกว้างกว่าคำศัพท์
2. ความสัมพันธ์ของหมวดหมู่กับคำศัพท์ต้องสูงมากพอเกินกว่าค่าหนึ่ง โดยที่ Term ทั้งหมดที่ได้จะ ถูกแบ่งออกเป็น 4 ชนิด

Type 0: หมวดหมู่และคำศัพท์เป็นคำๆเดียวกัน

Type 1: คำศัพท์มีความเกี่ยวข้องกับหมวดหมู่หมวดหมู่

Type 2: คำศัพท์อาจจะเกี่ยวข้องหรือไม่เกี่ยวข้องกับหมวดหมู่

Type 3: คำศัพท์และหมวดหมู่ไม่มีความเกี่ยวข้องกันเลย

เหตุผลที่เราต้องกำหนดชนิดของความสัมพันธ์เพราะว่าการหาว่า Relation degree มีมากแค่ไหน ในที่นี้เราจะใช้วิธีการ search เพื่อการหาความน่าจะเป็น และค่าความน่าจะเป็นที่เราจะ ได้สามารถนำมาจำแนกความสัมพันธ์ตามหมวดหมู่ทั้ง 4 เพื่อดูประสิทธิภาพของวิธีการทำ Relation test นี้ด้วย โดยวิธีการหาค่าความน่าจะเป็นสามารถหาได้ด้วยวิธีดังต่อไปนี้

เราสามารถหาค่า Probability จากการทำ Relation test จากสูตรต่อไปนี้

$$P(s|x) = H(s^x)/H(x) \quad (2.1)$$

โดยที่

s คือ หมวดหมู่ย่อย

เอกสาร x คือ คำศัพท์ภาษาอังกฤษ
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$P(s|x)$ คือ ค่าความน่าจะเป็นที่ s และ x จะมีความเกี่ยวข้องกัน

$H(x)$ คือ จำนวน hits ที่มี x ปรากฏอยู่

$H(s^x)$ คือ จำนวน hits ที่มีทั้ง s และ x ปรากฏอยู่

สมมุติฐานของการกำหนดค่าความน่าจะเป็นเพื่อให้เกี่ยวข้องกับหมวดหมู่คือ ถ้า $P(s|x) \geq 0.05$ หรือ $P(x|s) \geq 0.05$ แล้วคำศัพท์จะเกี่ยวข้องกับหมวดหมู่นี้อย่างใกล้ชิด สมมุติฐานเพิ่มเติมอื่นๆจะได้หลังจากการทำโครงการนี้เสร็จสิ้นและสรุปผลของการทำ relation test นี้



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 3

วิธีการดำเนินการศึกษา

ในการศึกษาโครงการนี้ เป็นการนำทฤษฎีการทำ Automatic Collection of Vocabulary and Related Term from the Web มาพัฒนาเป็น โปรแกรมคอมพิวเตอร์ เพื่อประยุกต์ใช้ในการทำ online dictionary รูปแบบใหม่ที่สามารถหาความเกี่ยวเนื่องกันของคำศัพท์ได้ด้วย โดยจะแบ่งการดำเนินการศึกษาออกเป็น 4 ส่วน ดังนี้

- โครงสร้างการทำงานตามทฤษฎี
- ขั้นตอนการดำเนินงาน
- การออกแบบโปรแกรมการทำงาน
- อัลกอริทึมในการทำงานของโปรแกรมตามทฤษฎี

3.1 โครงสร้างการทำงานตามทฤษฎี

โครงสร้างการทำงานตามทฤษฎี Automatic Collection of Vocabulary and Related Term from the Web ที่ใช้ในโครงการนี้ จะประกอบด้วยโครงสร้างการทำงานที่แบ่งเป็นจำนวนขั้นตอนการทำงานทั้งหมด 4 ขั้นตอน ดังนี้

1. การทำ Compiling Corpus
2. การทำ Automatic Term Recognition
3. การทำ Filtering
4. การทำ Relation test

3.2 ขั้นตอนการดำเนินงาน

โดยเริ่มจากกระบวนการของการวิเคราะห์การค้นหาคำศัพท์ จนกระทั่งถึงการหาความสัมพันธ์ของคำศัพท์และจัดเก็บลงฐานข้อมูล ซึ่งประกอบการทำงานได้เป็น 9 ขั้นตอนหลักๆ ดังนี้

1. กำหนดวัตถุประสงค์ (Objectives Determination) มีการศึกษาความต้องการในการวิเคราะห์ คือกำหนดปัญหาและวัตถุประสงค์การทำให้ชัดเจน ต้องเข้าใจถึงปัญหาและความต้องการของผู้ใช้งาน Online Dictionary ประกอบด้วยการวิเคราะห์เว็บไซต์เบื้องต้นว่ามีข้อมูลอะไรบ้างบนหน้าเว็บไซต์และเราต้องการคำศัพท์อะไรบ้างจากหน้าเว็บไซต์และเป็นคำศัพท์ที่จะนำมาใช้งานเมื่อใด
2. การจัดเตรียมรูปแบบของคำศัพท์ (Vocabulary preparation) ที่เราจะทำการค้นหาจากหน้า

เอกสารนี้เป็นเอกสารที่จัดทำขึ้นเพื่อใช้ในการศึกษาวิจัยเท่านั้น ไม่สามารถนำข้อมูลไปใช้
เอกสารนี้เป็นเอกสารที่จัดทำขึ้นเพื่อใช้ในการศึกษาวิจัยเท่านั้น ไม่สามารถนำข้อมูลไปใช้
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากเว็บไซต์ใดบ้าง มีเว็บไซต์ใดที่เหมาะสมและอยู่ในประเด็นที่ต้องการประกอบด้วย 2 ขั้นตอนย่อยคือ

2.1 การคัดเลือกคำศัพท์ที่ต้องการ (Vocabulary Selection) เป็นการกำหนดรูปแบบของคำศัพท์ว่าจะเอาคำศัพท์ประเภทใดบ้าง เป็นการเริ่มต้นของการค้นหาคำศัพท์ การเลือกคำศัพท์ขึ้นอยู่กับจุดประสงค์ของการใช้งานของเรา จึงเลือกคำศัพท์ทุกรูปแบบที่เจอและมีความหมายที่ถูกต้องแม่นยำ คำศัพท์ที่จะนำมาเก็บคือคำศัพท์ที่มีความหมายประกอบด้วยเท่านั้น

2.2 การตรวจสอบความถูกต้องของคำศัพท์และความหมายของคำศัพท์ที่ได้ (Vocabulary Preprocess) เป็นวิธีที่นำไปตรวจสอบคุณภาพของคำศัพท์และความหมายของคำศัพท์ที่ได้ โดยใช้วิธีการหลากหลายอาทิเช่น การนำคำศัพท์ไปค้นหาความหมายอีกครั้งว่าตรงกันกับความหมายที่ได้ในตอนแรกหรือไม่

3. การจัดเตรียมรูปแบบของการค้นหาคำศัพท์ (Vocabulary Platform) การวิเคราะห์วิธีการเลือกคำศัพท์ที่เหมาะสมกับแต่ ละเว็บไซต์ขั้นตอนนี้มีความสัมพันธ์กับการวิเคราะห์ข้อมูลในขั้นตอนที่ 1 ผ่านมา โดยอาจจะไปย้อนทำในขั้นตอนที่ 2 ใหม่ เนื่องจากถ้าวิธีการค้นหาคำศัพท์ได้ผลลัพธ์ไม่ตรงกับที่ต้องการก็ย้อนกลับไปทำใหม่จนกว่าจะได้คำศัพท์และความหมายที่ถูกต้อง
4. วิเคราะห์คำศัพท์ที่ได้ (Analysis of Result) การตรวจสอบคำศัพท์และคำแปลว่ามีความถูกต้องแม่นยำแค่ไหน เพื่อดูว่ารูปแบบของการค้นหาคำศัพท์ที่ได้มีประสิทธิภาพมากเพียงใด และประเมินผลลัพธ์ที่ได้จากขั้นตอนที่ 3
5. นำคำศัพท์ที่ได้มารวบรวมและจัดเก็บเข้าสู่ฐานข้อมูล โดยเรียงตามตัวอักษร A-Z
6. จัดทำหมวดหมู่ของคำศัพท์ โดยในรายงานชุดนี้จะมีการจัดทำหมวดหมู่ของคำศัพท์แบ่งเป็น 6 หมวดหมู่ใหญ่
7. หาความสัมพันธ์ของคำศัพท์และหมวดหมู่ที่จัดทำขึ้น ว่าคำศัพท์มีความสัมพันธ์กับหมวดหมู่ใด โดยมีความน่าจะเป็นของการอยู่ในหมวดหมู่นี้แค่ไหน เพื่อใช้ในการทำเป็นฐานความรู้ต่อไป
8. จัดเก็บคำศัพท์และคำแปล โดยเรียงตามหมวดหมู่ลงในฐานข้อมูล พร้อมเก็บค่าความน่าจะเป็นของความสัมพันธ์ของคำศัพท์ในแต่ละหมวดหมู่
9. วิเคราะห์ค่าความน่าจะเป็นได้ที่ได้ว่ามีความถูกต้องแม่นยำแค่ไหน เพื่อให้ในการปรับปรุงต่อไปในอนาคต
10. จัดทำเว็บไซต์ เพื่อให้บริการ Online Dictionary โดยสามารถค้นหาคำศัพท์แบบมาตรฐาน และการค้นหาคำศัพท์ที่ขึ้นต้นด้วยตัวอักษรที่ต้องการ, การค้นหาคำศัพท์ที่ประกอบด้วย

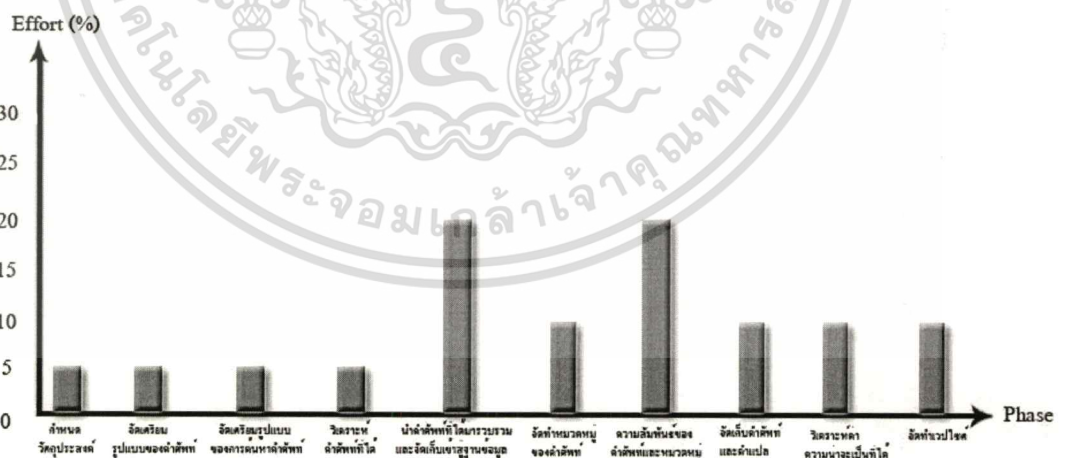
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตัวอักษรที่ต้องการ, การค้นหาคำศัพท์เรียงตามหมวดหมู่, การค้นหาคำศัพท์ที่ต้องการ คำศัพท์ที่อยู่ในหมวดหมู่เดียวกัน

การทำงานในแต่ละขั้นตอนจะใช้ระยะเวลาแตกต่างกันออกไป โดยจะใช้เวลาในการทำ ขั้นตอนต่างๆดังนี้

- การ กำหนดวัตถุประสงค์ 5%
- การจัดเตรียมรูปแบบของคำศัพท์ 5%
- การจัดเตรียมรูปแบบของการค้นหาคำศัพท์ 5%
- วิเคราะห์คำศัพท์ที่ได้ 5%
- นำคำศัพท์ที่ได้มารวบรวมและจัดเก็บเข้าสู่ฐานข้อมูล 20%
- จัดทำหมวดหมู่ของคำศัพท์ 10%
- หาความสัมพันธ์ของคำศัพท์และหมวดหมู่ที่จัดทำขึ้น 20%
- จัดเก็บคำศัพท์และคำแปลโดยเรียงตามหมวดหมู่ลงในฐานข้อมูล 10%
- วิเคราะห์ค่าความน่าจะเป็นได้ที่ได้ว่ามีความถูกต้องแม่นยำแค่ไหน 10%
- จัดทำเว็บไซต์ เพื่อให้บริการ Online Dictionary 10%

โดยสามารถแสดงเป็นกราฟดังภาพ 3.1 ดังนี้



รูปที่ 3.1 เปอร์เซ็นต์ของเวลาที่ใช้แต่ละขั้นตอน

การศึกษาโครงการนี้จะดำเนินการศึกษา Automatic Collection of Vocabulary and Related Term from the Web และทำการพัฒนาระบบสร้างเป็นโปรแกรมคอมพิวเตอร์เพื่อประยุกต์ใช้ในการสร้าง online dictionary ที่สามารถค้นหาคำศัพท์ที่มีความเกี่ยวข้องกัน และคำศัพท์ที่เกี่ยวข้องกับ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

และคำแปลแบบ automatic เพื่อความสะดวกและรวดเร็ว ในโครงการงานการนี้ การทำการตัดคำศัพท์ โดยการพัฒนาโปรแกรมขึ้นมา และมีการแยกโปรแกรมไว้เป็นส่วนๆเพื่อทำขั้นตอนต่างๆตาม ทฤษฎี Automatic Collection of Vocabulary and Related Term from the Web

5. การออกแบบฐานข้อมูล และ จัดเก็บคำศัพท์และคำแปลที่ได้ลงฐานข้อมูล ในโครงการงานนี้ เราจะมี การออกแบบฐานข้อมูล โดยจัดเก็บทั้งหมด 5 ตาราง ได้แก่ ตาราง Dictionary, ตาราง Directory, ตาราง Subdirectory, ตาราง Relation, ตาราง Link โดยที่แต่ละตารางมีการจัดเก็บข้อมูลดังต่อไปนี้

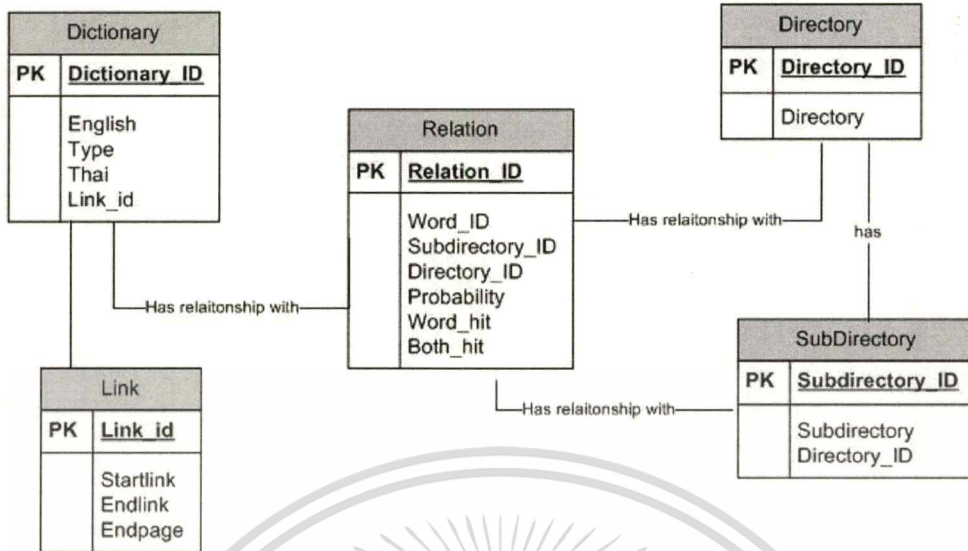
- ตาราง Dictionary จัดเก็บข้อมูลทั้งหมด 4 ข้อมูล ได้แก่ ID สำหรับเก็บรหัสของคำศัพท์, English สำหรับเก็บคำศัพท์, type สำหรับเก็บชนิดของคำศัพท์ และ Thai สำหรับเก็บคำแปลของ คำศัพท์

- ตาราง Directory จัดเก็บข้อมูลทั้งหมด 2 ข้อมูล ได้แก่ Directory_ID สำหรับเก็บรหัสของ หมวดหมู่, Directory สำหรับเก็บชื่อหมวดหมู่

- ตาราง Subdirectory จัดเก็บข้อมูลทั้งหมด 3 ข้อมูล ได้แก่ Subdirectory_ID สำหรับเก็บ รหัสของหมวดหมู่ย่อย, Subdirectory สำหรับเก็บชื่อหมวดหมู่ย่อย, Directory_ID สำหรับเก็บรหัส หมวดหมู่

- ตาราง Relation จัดเก็บข้อมูลทั้งหมด 7 ข้อมูล ได้แก่ Relation_ID สำหรับเก็บรหัสของ ความสัมพันธ์, Word_ID สำหรับเก็บรหัสของคำศัพท์ Subdirectory สำหรับเก็บชื่อหมวดหมู่ย่อย, Directory_ID สำหรับเก็บรหัสหมวดหมู่, Subdirectory_ID สำหรับเก็บรหัสของหมวดหมู่ย่อย, Directory_ID สำหรับเก็บรหัสของหมวดหมู่, Word_hit เก็บจำนวนที่เจอคำศัพท์, Both_hit เก็บ จำนวนที่เจอทั้งคำศัพท์และหมวดหมู่, Probability เก็บค่าความน่าจะเป็นที่คำศัพท์และหมวดหมู่จะ เกี่ยวข้องกัน

- ตาราง Link จัดเก็บข้อมูลทั้งหมด 3 ข้อมูล ได้แก่ Start_link , End_link และ End_page เพื่อ เก็บ URL ที่นำมาเปิดหา source file เพื่อใช้ในการค้นหาคำศัพท์ และความหมายของคำศัพท์ รูปแบบของฐานข้อมูล (Database diagram) สามารถดังแสดงได้ดังรูปที่ 3.5



รูปที่ 3.5 แสดงรูปแบบของฐานข้อมูล (Database diagram)

6. การคัดเลือกหมวดหมู่และหมวดหมู่ย่อย มาใช้งานในการจัดหมวดหมู่ของคำศัพท์ โดยในที่นี้ เราจะกำหนดให้มีหมวดหมู่ทั้งหมด 5 หมวดหมู่ ดังนี้

- Business & Finance คือ หมวดหมู่ที่เกี่ยวข้องกับ ธุรกิจและการเงิน ที่จะประกอบไปด้วย 20 หมวดย่อย
- Computers & Internet คือ หมวดหมู่ที่เกี่ยวข้องกับ คอมพิวเตอร์และอินเทอร์เน็ต ที่จะประกอบไปด้วย 15 หมวดย่อย
- Schools & Education คือ หมวดหมู่ที่เกี่ยวข้องกับ โรงเรียนและการศึกษา ที่จะประกอบไปด้วย 24 หมวดย่อย
- Entertainment & Arts คือ หมวดหมู่ที่เกี่ยวข้องกับบันเทิงและศิลปะ ที่จะประกอบไปด้วย 28 หมวดย่อย
- Health & Wellness คือ หมวดหมู่ที่เกี่ยวข้องกับ ธุรกิจและการเงิน ที่จะประกอบไปด้วย 19 หมวดย่อย
- Recreation & Sport คือ หมวดหมู่ที่เกี่ยวข้องกับ การพักผ่อนหย่อนใจและกีฬา ที่จะประกอบไปด้วย 9 หมวดย่อย

ในที่นี้เราจะระบุหมวดหมู่ในการทำในโครงการนี้ออกเป็น 6 หมวดหมู่ใหญ่ ดังต่อไปนี้

Business & Finance	Entertainment & Arts
Business Schools	Actors

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

<ul style="list-style-type: none"> • Companies • Consumer • Employment • Work • Financial Professionals • Free Stuff • Home Business • Industry Associations • Investments • Labor • Marketing • Advertising • News • Media • Public Relations • Real Estate • Small Business • Trade • Transportation 	<ul style="list-style-type: none"> • Actresses • Amusement • Theme Parks • Audio • Visual Equipment • Celebrities • Comics • Animation • Contests • Sweepstakes • Cool Clubs • Fashion • Fine Arts • Genres • Humanities • Humor • Magic • Movies • Official Groups • Clubs • Performing Arts • Quotations • Radio • Television • Trivia • Dancing • Music
Computers & Internet	
<ul style="list-style-type: none"> • Communications • Networking • Cyber culture • Data Formats • Desktop Publishing • Education • Hardware • Internet • Multimedia • Programming Languages • Security • Software • Technical Support • Collecting • Desktop Customization 	
Schools & Education	
<ul style="list-style-type: none"> • Bilingual • Subject • Classmates • Colleges • Universities • Courses • Distance Learning • Exchange Students 	
	Health & Wellness
	<ul style="list-style-type: none"> • Advice • Alternative Medicine • Beauty • Children's • Drugs • First Aid • Fitness • Health Care • Men's • Professional • Reproductive • Seniors • Stress Management • Support • Teens • Women's • Pet Health • Nutrition

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

<ul style="list-style-type: none"> · Fund Raising · Graduating Classes · Homework Help · Issues · Preschools · Religious Schools · Special Education · Student Organizations · Study Abroad · Teaching · Research · Testing · Theory · Methods · University Life · Dental Schools 	<ul style="list-style-type: none"> · Medications <p>Recreation & Sports</p> <ul style="list-style-type: none"> · Automotive · Aviation · Gambling · Outdoors · Sports · Toys · Travel · Disabled · Games
---	--

โดยจะนำมาจัดเก็บในฐานข้อมูลโดยแบ่งเป็น 2 ตาราง ได้แก่ตาราง Directory และ ตาราง Subdirectory โดยที่จะมีการให้ค่า Directory_id ไว้ในตาราง Subdirectory เพื่อบ่งบอกว่าเป็นหมวดหมู่ย่อยของหมวดหมู่ใด ดังตัวอย่างรูป 3.6

Directory : Table		Subdirectory : Table		
Directory_id	Directory	Subdirectory_id	Subdirectory	Directory_id
1	Business & Finance	1	Communications	2
2	Computers & Internet	2	Cyber culture	2
3	Schools & Education	3	Data Formats	2
4	Entertainment & Arts	4	Desktop Publishing	2
5	Health & Wellness	5	Education	2
6	Recreation & Sports	6	Hardware	2
		7	Internet	2
		8	Multimedia	2
		9	Programming Languages	2
		10	Security	2
		11	Software	2
		12	Technical Support	2
		13	Collecting	2
		14	Desktop Customization	2
		15	Bilingual	3
		16	By Subject	3
		17	Classmates	3
		18	Colleges	3
		19	Courses	3
		20	Distance Learning	3
		21	Exchange Students	3
		22	Fund Raising	3
		23	Graduating Classes	3

รูปที่ 3.6 แสดงวิธีการเก็บหมวดหมู่และหมวดหมู่ย่อยลงในฐานข้อมูล

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

7. การทำ Relation test และการเก็บผลลัพธ์ที่ได้ เราจะทำการหาค่าทั้งหมด 3 ค่า คือค่า Word hits และค่า both hits เพื่อนำมาหาค่า Probability ที่เราต้องการ โดยเราสามารถคำนวณหาแต่ละค่าได้ด้วยวิธีการดังต่อไปนี้

- ค่า Word hits คือการนำคำศัพท์ที่เราต้องการ ไปหาค่า hits โดยใช้ search engine hits โดยที่เราจะนำคำศัพท์ไป search ในเว็บ search engine และนับจำนวน hits ที่เจอโดยเลือกเฉพาะ 30 ลิงค์แรกที่เจอ และนับเฉพาะ ลิงค์ที่มีคำศัพท์ที่เราต้องการปรากฏอยู่ในส่วนของลิงค์เท่านั้น เราจะไม่นับถ้าไม่เจอในส่วนของลิงค์ ถึงแม้ว่าจะเจอในส่วนของ content ก็ตาม ซึ่งในโครงงานนี้ เราจะพัฒนาโปรแกรมอันโนมัติขึ้นมาทำการคำนวณค่า Word hits ตามวิธีดังกล่าว ดังตัวอย่างรูป 3.7

Web Images Groups News more »

Google™ Fitness Search Advanced Search Preferences

Search: the web pages from Singapore

Web

Fitness Equipment, Workout, **Fitness** Program, **Fitness** Articles ...
Fitness Equipment, Workout, **Fitness**, **Fitness** Forum, **Fitness** Program, **LA Fitness**,
Fitness Articles, Beauty **Fitness**, Weight Loss, Weight Training.
www.fitness.com/ - 35k - Cached - Similar pages

Fitness Online
Fit Pregnancy presents our very own postnatal **fitness** DVD. ... Men's **Fitness**
Muscle Prescription It's the perfect plan for building your perfect body - fast ...
www.fitnessonline.com/ - 34k - 6 Feb 2007 - Cached - Similar pages

Fitness Magazine Home Page
Fitness, beauty, wellness, food and other topics. Includes recipe finder and email
newsletter.
www.fitnessmagazine.com/ - 53k - 6 Feb 2007 - Cached - Similar pages

The President's Council on Physical **Fitness** and Sports - You're it ...
Information to encourage Americans to become physically active and participate in
sports. Includes details of the President's Challenge award programs for ...
www.fitness.gov/ - 27k - Cached - Similar pages

LA Fitness: Where **Fitness** is a Way of Life
At **LA Fitness** we want you to exercise your options. From basketball to racquetball,
swimming to indoor cycling, free weights to cardio, personal training to ...
www.lafitness.com/ - 10k - Cached - Similar pages

Health and **fitness** clubs, gym membership: **LA fitness**
Health and **fitness** clubs offering gym membership, personal **fitness** training,
corporate gym membership and exercise classes.
www.lafitness.co.uk/ - 18k - Cached - Similar pages

Fitness in the Yahoo! Directory
Features sites that focus on physical **fitness**, good nutrition, and exercise. Includes
sites for health clubs, **fitness** experts, magazines, equipment and ...
dir.yahoo.com/Health/**Fitness**/ - 15k - Cached - Similar pages

Open Directory - Health **Fitness**
the entire directory, only in Health/**Fitness**. Top: Health: **Fitness** (1027) ... Usenet
misc.**fitness**.misc - news: - Google Groups. "**Fitness**" search on: ...
dmoz.org/Health/**Fitness**/ - 10k - Cached - Similar pages

รูปที่ 3.7 แสดงวิธีการนับจำนวน word hits จาก search engine

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- ค่า Both hits คือการนำคำศัพท์ และหมวดหมู่ที่เราต้องการ ไปหาค่า hits โดยใช้ search engine hits โดยที่เราจะนำคำศัพท์และหมวดหมู่ไป search ในเว็บ search engine และนับจำนวน hits ที่เจอโดยเลือกเฉพาะ 30 ลิงค์แรกที่เจอ และนับเฉพาะ ลิงค์ที่มีคำศัพท์และหมวดหมู่ที่เราต้องการ ปรากฏอยู่ในส่วนของลิงค์เท่านั้น เราจะไม่นับถ้าไม่เจอในส่วนของลิงค์ ถึงแม้ว่าจะเจอในส่วนของ content ก็ตาม โดยที่รูปแบบของการค้นหาคือ “คำศัพท์” “หมวดหมู่” ซึ่งในโครงการนี้ เราจะพัฒนา โปรแกรมอัน โนมัติขึ้นมาทำการคำนวณค่า Word hits ตามวิธีดังกล่าว ดังตัวอย่างรูป3.8

The screenshot shows a Google search interface with the search term 'Fitness tennis'. Below the search bar, there are several search results listed under the 'Web' category. Each result includes a title, a brief description, and the URL. The results are as follows:

- Tennis training program to improve your tennis fitness**
Training program to help improve your racket game, all for FREE from netfit.
www.netfit.co.uk/racquets-web.htm - 47k - Cached - Similar pages
- Tennis Fitness Testing**
A discussion of the fitness tests required for the sport of tennis.
www.topendsports.com/sport/tennis/testing.htm - 19k - Cached - Similar pages
- Fitness and Tennis Club, Norwalk, CT**
Since 1979, Fitness and Tennis Club has been providing services as a Full Service Health Club in the Norwalk, CT area.
www.fitnessandtennisclub.com/ - 22k - Cached - Similar pages
- Amazon: Listmania! - View List "Top 5 Tennis Fitness Books"**
Net Flex: 10 Minutes a Day to Better Play by Paul Frediani. \$9.95 Used & New from: \$5.09. Average Customer Rating: Good Tennis fitness book. ...
www.amazon.com/Top-5-Tennis-Fitness-Books/lm/1VHL6JZW4XD5 - 61k - Cached - Similar pages
- Baltimore Fitness and Tennis --- FitnessandTennis.com | Home**
Baltimore Fitness & Tennis is a modern, state of the art fitness and tennis facility that offers fitness and tennis memberships, innovative fitness and ...
fitnessandtennis.com/ - 21k - Cached - Similar pages
- Libertville Tennis and Fitness**
A quality tennis and fitness club, with a long tradition of providing the winning edge in personalized programs and world class instruction.
www.ltf.com/ - 38k - Cached - Similar pages
- New Hampshire - NH Health Club, NH Fitness Center and NH Tennis ...**
New Hampshire health club, fitness center and tennis club providing fun filled activities to New Hampshire residents and others.
www.mountainsiderfc.com/ - 21k - Cached - Similar pages
- Sun Oaks Tennis and Fitness Online**
Sun Oaks Tennis & Fitness is located at 3452 Argyle Rd. Redding, CA 96002. You may reach us by calling (530) 221-4405 or by email at info@sunoaks.com. ...
☎ Map of 3452 Argyle Rd, Redding, CA 96002, USA
www.sunoaks.com/ - 12k - Cached - Similar pages

รูปที่3.8 แสดงวิธีการนับจำนวน both hits จาก search engine

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

-ค่า Probability คือค่าความน่าจะเป็นที่คำศัพท์และหมวดหมู่จะมีความเกี่ยวข้องกัน โดยที่เราจะแยกประเภทของความสัมพันธ์ออกเป็น 4 ประเภท โดยเราจะกำหนดสมมุติฐานของการทำลองในโครงการนี้ โดยแบ่งตามค่า probability ดังนี้

Probability =1

อยู่ใน Type 0: หมวดหมู่และคำศัพท์เป็นค่าๆเดียวกัน หรือมีความเกี่ยวข้องกันมาก

0.8=< Probability <1

อยู่ใน Type 1: คำศัพท์มีความเกี่ยวข้องกับหมวดหมู่หมวดหมู่

0.5=<Probability <0.8

อยู่ใน Type 2: คำศัพท์อาจจะเกี่ยวข้องหรือไม่เกี่ยวข้องกันกับหมวดหมู่

Probability <0.5

อยู่ใน Type 3: คำศัพท์และหมวดหมู่ไม่มีความเกี่ยวข้องกันเลย

โดยที่เราจะแบ่งการพิจารณาออกจากรูปแบบของคำศัพท์ คือคำนาม คำกริยา และคำอื่นๆ เนื่องจากคำที่จะถูกต้องมากที่สุดในการนำมาทำ Relation test ตามสมมุติฐานของโครงการนี้คือ คำนาม ดังนั้นเราจะพิจารณาเฉพาะคำนามเป็นหลัก แต่จะมีคำชนิดอื่นๆ ในการวิเคราะห์ผลลัพธ์ในโครงการนี้ด้วย ซึ่งในส่วนของการทดลองครั้งนี้จะแบ่งผลการทำ Relation ออกเป็นคำนาม กับคำชนิดอื่นๆ เพื่อใช้แสดงในการค้นหาผ่านเว็บไซต์ online dictionary และเพื่อใช้ในการวิเคราะห์ความถูกต้องของทฤษฎีการทำ Automatic Collection of Vocabulary and Related Term from the Web

เราสามารถหาค่า Probability จากการทำ Relation test จากสูตรต่อไปนี้

$$P(s|x) = H(s^x)/H(x) \quad (3.1)$$

โดยที่

s คือ หมวดหมู่ย่อย

x คือ คำศัพท์ภาษาอังกฤษ

$P(s|x)$ คือ ค่าความน่าจะเป็นที่ s และ x จะมีความเกี่ยวข้องกัน

$H(s)$ คือ จำนวน hits ที่มี s ปรากฏอยู่

$H(x)$ คือ จำนวน hits ที่มี x ปรากฏอยู่

$H(s^x)$ คือ จำนวน hits ที่มีทั้ง s และ x ปรากฏอยู่

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

โดยจะนำค่าทั้งหมดที่ได้ ไปเก็บลงฐานข้อมูลดังตัวอย่างรูป 3.9

Relation : Table						
relation_id	Word_id	Subdirectory_id	Directory_id	probability	Word_hit	Both_hit
20780	20228	45	1	0.03703703703	27	1
20781	20228	46	1	0	27	0
20782	20228	47	1	0	27	0
20783	20228	48	1	0	27	0
20784	20228	49	1	0	27	0
20785	20228	50	1	0.07407407407	27	2
20786	20228	51	1	0.07407407407	27	2
20787	20228	52	1	0	27	0
20788	20229	36	1	0	27	0
20789	20229	37	1	0.59259259259	27	16
20790	20229	38	1	0.62962962962	27	17
20791	20229	39	1	0.03703703703	27	1
20792	20229	40	1	0.40740740740	27	11
20793	20229	41	1	0.03703703703	27	1
20794	20229	42	1	0.14814814814	27	4
20795	20229	43	1	0.03703703703	27	1
20796	20229	44	1	0.44444444444	27	12
20797	20229	45	1	0.59259259259	27	16
20798	20229	46	1	0	27	0
20799	20229	47	1	0.22222222222	27	6
20800	20229	48	1	0.29629629629	27	8
20801	20229	49	1	0.29629629629	27	8
20802	20229	50	1	0.11111111111	27	3
20803	20229	51	1	0.62962962962	27	17
20804	20229	52	1	0.59259259259	27	16

รูปที่ 3.9 แสดงการเก็บค่าต่างๆที่ได้จากการทำ Relation test

จากรูป 3.9 จะเห็นได้ว่าเราจะได้อ่านค่า Probability ต่างๆกันสำหรับคำศัพท์และหมวดหมู่ต่างๆ ซึ่งจะนำค่าเหล่านี้ไปวิเคราะห์ในส่วนของผลการทดลอง เพื่อใช้ในการสรุปสมมติฐานของค่า Probability สำหรับพัฒนาการทำ Relation test ต่อไปในอนาคต

7. การวิเคราะห์ผลลัพธ์ที่ได้จากการทำ Relation test เราจะนำฐานข้อมูลที่ได้มาทำการวิเคราะห์ผล ด้วยการวิเคราะห์แบบ manual check ว่าคำศัพท์และหมวดหมู่ที่ได้ มีผลตามสมมติฐานที่เราตั้งไว้หรือไม่ และควรปรับสมมติฐานอย่างไร เพื่อให้เหมาะสมและได้ผลลัพธ์ที่ถูกต้องที่สุด พร้อมทั้งยังหาค่าของความถูกต้องของวิธีการทำ Relation test นี้ว่ามีความถูกต้องมากน้อยเพียงใด โดยการนำค่า Probability ที่ได้จากแต่ละแถว มาพิจารณาว่าคำศัพท์ และ หมวดหมู่อะไร มีความสัมพันธ์กันอย่างไร ในแต่ละช่วงค่าของ Probability ที่ได้

8. การสร้าง Online dictionary เพื่อนำฐานข้อมูลที่ได้มาใช้ให้เกิดประโยชน์สูงสุด โดยจะจัดทำ online dictionary ที่สามารถค้นหาคำศัพท์ได้ 3 ลักษณะ ได้แก่

- แบบมาตรฐาน ได้แก่การค้นหาคำจากส่วนเริ่มต้นของคำ, การค้นหาคำจากส่วนประกอบของ

คำ, การค้นหาคำจากส่วนท้ายสุดของคำ, ค้นหาเฉพาะเจาะจงคำนั้นๆเท่านั้น

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์ไว้เพื่อใช้ในการศึกษาวิจัยเท่านั้น ไม่ให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- แบบการค้นหาคำศัพท์ที่ relation กัน โดยเราจะค้นหาคำศัพท์ในหมวดหมู่เดียวกันขึ้นมาแสดงทั้งหมด เพื่อช่วยให้ผู้ใช้งานนำคำศัพท์เหล่านี้ไปใช้งานในการเขียนหรือการอ่านได้ง่ายมากยิ่งขึ้น

- แบบการค้นหาคำศัพท์ที่เกี่ยวข้องกับหมวดหมู่ โดยเลือกหมวดหมู่ที่มีจัดเตรียมไว้ให้เพื่อดูคำศัพท์ที่มีความเกี่ยวข้องกัน เพื่อใช้ประโยชน์ในด้านต่างๆ

ดังแสดงขั้นตอนการดำเนินงานของโปรแกรมได้ดังตัวอย่างรูป 3.10



รูปที่3.10 แสดงขั้นตอนการดำเนินงาน

3.3 การพัฒนาโครงการตามลำดับขั้น

การพัฒนาโครงการนี้ได้แบ่งการทำงานเป็น 3 ลำดับขั้น ซึ่งพัฒนาตามผลการทดลองของโครงการเพื่อให้ผลลัพธ์ที่มีความถูกต้องสูงขึ้น เอกสารนี้เป็นเอกสารที่ลงนามแล้วซึ่งใช้ในการศึกษาเท่านั้น ไม่นอญญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เท่าที่ผู้ต้องการ, ไม่มีขอบเขต, เดินเฟ้นฟ่าน, จังหะเต็นรำอะ โก โก่อ้, เวทีเต็นรำจัังหะอะ โก โโก. - S...

จะทำการตัดคำศัพท์โดยตัดคำที่อยู่ด้านหน้า มาเป็นคำศัพท์ และ หลังจาก มาเป็นคำแปลของคำศัพท์แต่ละคำ

5. เก็บรวบรวมคำศัพท์และคำแปลลงฐานข้อมูล

6. เลือกหมวดหมู่หลักมา 6 หมวดหมู่ และหมวดหมู่ย่อยมาจำนวน 115 หมวดหมู่โดยนำมาจากหมวดหมู่ของ www.yahoo.com

7. ทำ Relation test ตามสูตร $P(s|x) = H(s^x)/H(x)$ โดยใช้ค่า hits ที่ได้จากการค้นหาจาก www.google.com

8. เก็บผลลัพธ์จากการทำ Relation test ของทุกๆคำศัพท์กับทุกๆหมวดหมู่ย่อยที่เราได้

9. วิเคราะห์ความถูกต้องของผลลัพธ์โดยการทำ manual test โดยพิจารณาจากการใช้ search engine เข้าช่วยในการพิจารณาว่าค่า Relation ระหว่างผลลัพธ์และหมวดหมู่มีความเกี่ยวข้องกันมากน้อยเพียงใด โดยนำคำศัพท์และหมวดหมู่มาทำการค้นหาใน www.google.com โดยใส่เงื่อนไขในการค้นหาคือ “หมวดหมู่”+”คำศัพท์” และดูว่าพบมากน้อยเพียงใด และรวมถึงการตรวจสอบจากการพิจารณาแบบ manual test

จากการทำงานตามวิธีการข้างต้นจะสามารถสรุปผลลัพธ์ได้ดังนี้

1. ในเวลา 11 นาทีสามารถหาความสัมพันธ์ของคำศัพท์ได้ 1 คำ กับ หมวดหมู่ย่อย 115 หมวดหมู่

2. ผลลัพธ์ที่ได้สำหรับคำศัพท์ชนิดต่างๆมีความแตกต่างกัน สำหรับคำศัพท์ที่เป็นคำนามจะมีความแม่นยำมากที่สุด คำสรรพนาม คำกริยา และคำชนิดอื่นๆจะมีความแม่นยำน้อยมาก

3. หมวดหมู่ย่อยจำนวนมากที่มีประโยชน์ต่อการใช้งานน้อย และไม่ชี้เฉพาะเจาะจงเพียงพอสำหรับการนำมาหาค่าความสัมพันธ์

4. การทำ Relation ทุกคำศัพท์ที่เก็บไว้กับทุกหมวดหมู่ใช้เวลามากและได้ผลลัพธ์ที่มีความแม่นยำน้อย

- การพัฒนาโครงการในลำดับที่ 2 คือการพัฒนาโครงการจากการพัฒนาโครงการในลำดับที่ 1 ให้มีความถูกต้องของผลลัพธ์มากขึ้น โดยมีการทำงานดังนี้

1. คัดเลือก online dictionary สำหรับการค้นหาคำศัพท์และคำแปลโดยการทำแบบ manual โดยการเลือกจากเว็บไซต์ที่มีคำศัพท์ปริมาณมากที่สุด และมีการแสดงคำศัพท์และคำแปลใน source ของหน้า website ที่สามารถเรียกขึ้นมาดูได้ และมีรูปแบบการแสดงผลที่ที่แน่นอน

2. รวบรวมคำศัพท์และเก็บรวบรวมเป็น corpus โดยใช้ฟังก์ชันการทำงาน

wc.DownloadFile(URL)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3. ตัดประโยคที่มีคำศัพท์และคำแปลโดยใช้การทำ sentence extraction โดยดูจากเงื่อนไขของภาษา html ตามแต่ละ online dictionary ที่เรานำมาใช้ โดยพิจารณาจาก source ใน corpus ที่เรารวบรวมได้ดังเช่นเดียวกับการพัฒนาโครงการในลำดับที่ 1

4. ตัดคำศัพท์และคำแปลโดยใช้การทำ sentence extraction โดยดูจากเงื่อนไขของภาษา html ตามแต่ละ online dictionary ที่เรานำมาใช้ดังเช่นเดียวกับการพัฒนาโครงการในลำดับที่ 1

5. เก็บรวบรวมคำศัพท์และคำแปลลงฐานข้อมูล

6. เลือกหมวดหมู่หลักมา 6 หมวดหมู่ และหมวดหมู่ย่อยมาจำนวน 50 หมวดหมู่จากเดิม 115 หมวดหมู่ โดยเลือกเฉพาะหมวดหมู่ที่มีความน่าสนใจและสามารถนำมาใช้ประโยชน์ได้

7. ทำการคัดกรองคำศัพท์เฉพาะคำนาม โดยดูจากคำแปลจะมีอักษร n. หรือ (n) เพื่อบอกว่าเป็นชนิดคำนาม โดยคัดเลือกมาทั้งหมด 1326 คำ โดยวิธี Random เพื่อให้ได้คำนามที่มีความหลากหลายและต่างชนิดกัน

8. ทำ Relation test ตามสูตร $P(s|x) = H(s^x)/H(x)$ โดยใช้ค่า hits ที่ได้จากการค้นหาจาก www.google.com

9. เก็บผลลัพธ์จากการทำ Relation test ของทุกๆ คำศัพท์กับทุกๆ หมวดหมู่ย่อยที่เราได้ผลลัพธ์ทั้งหมด 60000 Relation

10. วิเคราะห์ความถูกต้องของผลลัพธ์โดยการทำ manual test โดยพิจารณาจากการใช้ search engine เข้าช่วยในการพิจารณาว่าค่า Relation ระหว่างผลลัพธ์และหมวดหมู่มีความเกี่ยวข้องกันมากน้อยเพียงใด โดยนำคำศัพท์และหมวดหมู่มาทำการค้นหาใน www.google.com โดยใส่เงื่อนไขในการค้นหาคือ “หมวดหมู่”+”คำศัพท์” และดูว่าพบมากน้อยเพียงใด และรวมถึงการตรวจสอบจากการพิจารณาแบบ manual test

จากการทำงานตามวิธีการข้างต้นจะสามารถสรุปผลลัพธ์ได้ดังนี้

1. ในเวลา 3 นาทีสามารถหาความสัมพันธ์ของคำศัพท์ได้ 1 คำ กับ หมวดหมู่ย่อย 50 หมวดหมู่

2. ผลลัพธ์ที่ได้สำหรับคำศัพท์ชนิดคำนามมีความถูกต้องแม่นยำมากกว่าการพัฒนาในลำดับขั้นตอนแรก แต่ยังไม่ได้ผลลัพธ์ที่น่าพอใจ

3. การคัดเลือกหมวดหมู่ย่อยเพื่อนำมาใช้งานทำให้ช่วยประหยัดเวลาและได้ผลลัพธ์ที่แม่นยำขึ้น

4. การพัฒนาโครงการในลำดับที่ 2 ได้ผลลัพธ์ที่แม่นยำกว่าการพัฒนาโครงการในลำดับที่ 1 แต่ยังไม่มีความแม่นยำเท่าที่ควร โดยความแม่นยำที่ควรจะได้คือ 70% ขึ้นไป

- การพัฒนาโครงการในลำดับที่ 3 คือการพัฒนาโครงการจากการพัฒนาโครงการในลำดับที่ 2 ให้มีความถูกต้องของผลลัพธ์มากขึ้น โดยมีการทำงานดังนี้

เอกสารนี้เป็นเอกสารสงวนลิขสิทธิ์สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1. คัดเลือก online dictionary สำหรับการค้นหาคำศัพท์และคำแปลโดยการทำแบบ manual โดยการเลือกจากเว็บไซต์ที่มีคำศัพท์ปริมาณมากที่สุด และมีการแสดงคำศัพท์และคำแปลใน source ของหน้า website ที่สามารถเรียกขึ้นมาดูได้ และมีรูปแบบการแสดงผลที่ที่แน่นอน

2. รวบรวมคำศัพท์และเก็บรวบรวมเป็น corpus โดยใช้ฟังก์ชันการทำงาน
wc.DownloadFile(URL)

3. ตัดประโยคที่มีคำศัพท์และคำแปลโดยใช้การทำ sentence extraction โดยดูจากเงื่อนไขของภาษา html ตามแต่ละ online dictionary ที่เรานำมาใช้ โดยพิจารณาจาก source ใน corpus ที่เรารวบรวมได้ดังเช่นเดียวกับการพัฒนาโครงการในลำดับที่ 1 และ 2

4. ตัดคำศัพท์และคำแปลโดยใช้การทำ sentence extraction โดยดูจากเงื่อนไขของภาษา html ตามแต่ละ online dictionary ที่เรานำมาใช้ดังเช่นเดียวกับการพัฒนาโครงการในลำดับที่ 1 และ 2

5. เก็บรวบรวมคำศัพท์และคำแปลลงฐานข้อมูล

6. เลือกหมวดหมู่หลักมา 6 หมวดหมู่ และหมวดหมู่ย่อยมาจำนวน 50 หมวดหมู่จากเดิม 115 หมวดหมู่

7. ทำการคัดกรองคำศัพท์เฉพาะคำนาม โดยดูจากคำแปลจะมีอักษร n. หรือ (n) เพื่อบอกว่าเป็นชนิดคำนาม โดยคัดเลือกมาทั้งหมด 1326 คำ โดยวิธี Random เพื่อให้ได้คำนามที่มีความหลากหลายและต่างชนิดกัน

8. ทำ Relation test ตามสูตร $P(s|x) = H(s^x)/H(x)$ โดยเปลี่ยนแปลงจากการใช้ค่า hits เป็นการพิจารณาเฉพาะ 30 ผลลัพธ์แรกที่ได้ และใช้วิธีการนับจำนวนที่เจอคำศัพท์และหมวดหมู่ย่อยในหัวข้อของผลลัพธ์ที่เป็นค่า $H(s^x)$ และ ใช้วิธีการนับจำนวนที่เจอหมวดหมู่ย่อยในหัวข้อของผลลัพธ์เป็นค่า $H(s)$

9. เก็บผลลัพธ์จากการทำ Relation test ของทุกๆคำศัพท์กับทุกๆหมวดหมู่ย่อยที่เราได้ผลลัพธ์ทั้งหมด 60000 Relation

10. วิเคราะห์ความถูกต้องของผลลัพธ์โดยการทำ manual test โดยพิจารณาจากการใช้ search engine เข้าช่วยในการพิจารณาว่าค่า Relation ระหว่างผลลัพธ์และหมวดหมู่มีความเกี่ยวข้องกันมากน้อยเพียงใด โดยนำคำศัพท์และหมวดหมู่มาทำการค้นหาใน www.google.com โดยใส่เงื่อนไขในการค้นหาคือ “หมวดหมู่”+”คำศัพท์” และดูว่าพบมากน้อยเพียงใด และรวมถึงการตรวจสอบจากการพิจารณาแบบ manual test

จากการทำงานตามวิธีการข้างต้นจะสามารถสรุปผลลัพธ์ที่ได้ดังนี้

1. ในเวลา 3 นาทีสามารถหาความสัมพันธ์ของคำศัพท์ได้ 1 คำ กับ หมวดหมู่ย่อย 50 หมวดหมู่

2. ผลลัพธ์ที่ได้สำหรับคำศัพท์ชนิดคำนามมีความถูกต้องแม่นยำมากกว่าการพัฒนาในลำดับขั้นตอนที่ 2 และจาก 10 คำที่มีค่า Relation มากกว่า 8 มีได้ผลลัพธ์ที่แม่นยำ 70%

3. การพัฒนา โดยการเปลี่ยนจากการหาค่า hits มาเป็นการหาค่าที่พิจารณาเฉพาะ 30 ผลลัพธ์แรกที่ได้ ทำให้ได้ผลลัพธ์ที่แม่นยำและถูกต้องมากขึ้น

4. การพัฒนา โครงการงานในลำดับที่3 ได้ผลลัพธ์ที่แม่นยำกว่าการพัฒนาโครงการงานในลำดับที่1 และการพัฒนาโครงการงานในลำดับที่2 และแม่นยำเท่าที่ควรจึงพร้อมสำหรับนำไปทำ online dictionary เพื่อแสดงผลลัพธ์ของโครงการงานนี้

3.4 รูปแบบของผลลัพธ์ตามขั้นตอนการทำงาน

ในส่วนของ Application

- การทำ Compiling Corpus ผลลัพธ์ที่ได้คือชุดของข้อมูลทั้งหมดที่ได้จากการรวบรวมหน้าเว็บไซต์
- การทำ Automatic Term Recognition ผลลัพธ์ที่ได้คือชุดของข้อมูลที่มีขนาดเล็กลงที่ได้จากการตัดเฉพาะส่วนของข้อมูลที่มีคำศัพท์และคำแปลอยู่ข้างใน
- การทำ filtering ผลลัพธ์ที่ได้คือคำศัพท์และคำแปลสำหรับการเก็บลงฐานข้อมูล
- การทำ Relation test ผลลัพธ์ที่ได้คือค่า Relation และเก็บรวบรวมลงฐานข้อมูล

ในส่วนของ Website

- การทำ online dictionary เพื่อแสดงผลลัพธ์ในส่วนของ online dictionary มาตรฐาน คือ การค้นหาความหมายของคำศัพท์
- การทำ online dictionary เพื่อแสดงผลลัพธ์ในส่วนของ relation คือการแสดงว่าศัพท์แต่ละคำมีความเกี่ยวข้องกับหมวดหมู่ใดบ้าง

บทที่ 4

การทำงานของโครงการและผลการทำงาน

จากการทดลอง สามารถสรุปผลการทดลองการทำงานจากการพัฒนา โปรแกรม Automatic Collection of Vocabulary and Related Term from the Web โดยแบ่งการอธิบายออกเป็นส่วนต่างๆ ดังนี้

1. ส่วนของการเตรียมข้อมูล
2. ส่วนการทำงานของโปรแกรมในส่วนของ back office
3. ส่วนการทำงานของโปรแกรมในส่วนของเว็บไซต์
4. ส่วนของการแสดงผลลัพธ์

4.1 การเตรียมข้อมูล

ในการทดลองนี้ได้มีการเตรียมข้อมูลที่จะใช้ในการทดลองเพื่อการทำ Automatic Collection of Vocabulary and Related Term from the Web โดยมีการแบ่งการเตรียมข้อมูลเป็น 2 ส่วน คือส่วนแรกคือข้อมูลที่จัดเตรียมเพื่อมาทำเป็น online dictionary และส่วนที่สองคือข้อมูลที่เรารวบรวมเพื่อใช้ในการหา related term from website ของคำศัพท์และหมวดหมู่ที่เรากำหนด สิ่งที่เราจะต้องจัดเตรียมในขั้นตอนนี้คือการจัดเตรียมข้อมูลในส่วนของคำศัพท์และหมวดหมู่ ในที่นี้เรามีคำศัพท์ในฐานะข้อมูลเป็นจำนวนมากกว่า 70,000 คำ และมีหมวดหมู่ย่อยทั้งหมด 118 หมวดหมู่ จากผลการทดลองการหาค่าความสัมพันธ์ทำให้ทราบว่าจากทฤษฎีนี้สามารถหาค่าความสัมพันธ์ของหมวดหมู่และคำศัพท์ได้ถูกต้องในระดับที่ดี โดยการเตรียมข้อมูลสำหรับใช้ในการทดลองนี้เป็นสิ่งสำคัญต่อการทำงานของการทำงานของการหาค่า Relation test ดังนั้นคำศัพท์และหมวดหมู่ที่นำมาหาค่าความสัมพันธ์จะผ่านขั้นตอนในการจัดเตรียมให้เป็นข้อมูลที่พร้อมในการใช้งาน จากนั้นต้องนำไปผ่านกระบวนการคัดเลือกข้อมูล เพื่อให้ข้อมูลพร้อมใช้สำหรับ โปรแกรม Automatic Collection of Vocabulary and Related Term from the Web ในที่นี้มีการเตรียมข้อมูลสำหรับในการทำ Relation test เราจัดคัดเลือกคำศัพท์มาทำ เป็นจำนวน 1300 ตัวเฉพาะที่เป็นชนิดคำนาม และ หมวดหมู่หลัก 6 หมวดหมู่ และ หมวดหมู่ย่อย 46 หมวดหมู่ นำมาทำ relation test ซึ่งจะได้ข้อมูลทั้งหมดคือ $1300 * 20$ คือ 59800 relation เพื่อนำมาใช้ในการเปรียบเทียบดูผลการทดลอง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.1 แสดงข้อมูลที่ใช้ในการทดลองของ Directory

Directory : Table	
Directory_id	Directory
1	Business & Finance
2	Computers & Internet
3	Schools & Education
4	Entertainment & Arts
5	Health & Wellness
6	Recreation & Sports

ตารางที่ 4.2 แสดงข้อมูลที่ใช้ในการทดลองของ Subdirectory

Hardware	First Aid
Internet	Fitness
Programming	Health Care
Software	Pet Health
Technical	Automotive
Computer	Aviation
Subject	Gambling
Courses	Outdoor
Student	Sport
Graduate	Toy
Schools	Travel
Education	Game
University	Business
Actors	Company
Audio	Consumer
Comics	Employment
Fashion	Financial
Movies	Business
Radio	Industry
Television	Marketing
Dancing	Drug
Music	First Aid
Advice	Fitness
Beauty	Health Care
Children	Pet Health

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.3 แสดงข้อมูลที่ใช้ในการทดลองของ Dictionary

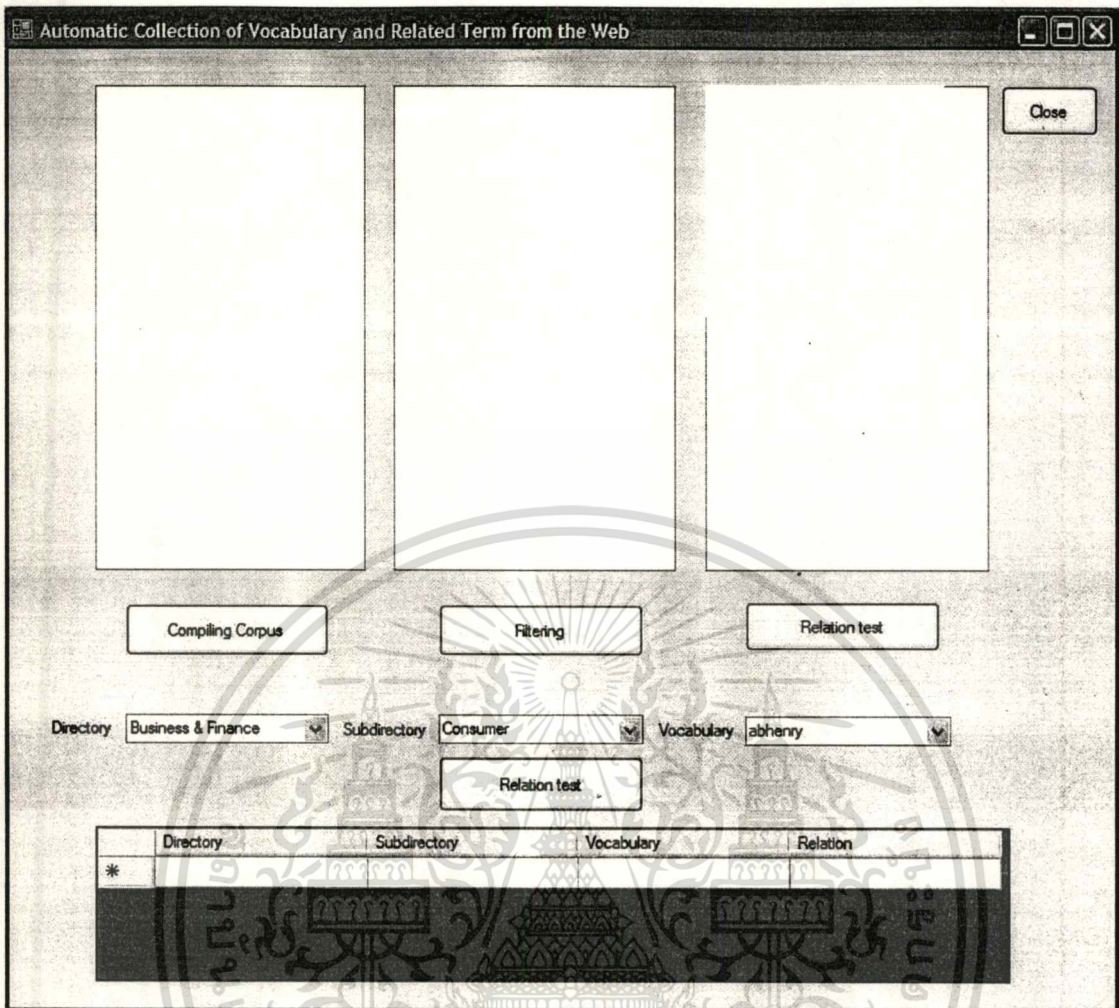
English	Thai
abhenry	(แอมเฮน'รี) น. หน่วยการนำไฟฟ้าที่เป็นเซนต์เมตร-กรัม-วินาที และมีค่าเท่ากับ 10-
ablegate	(แอม'ลิเดท) น. ผู้แทนของโบป
accusation	(แอคคิวเซ'ชัน) น. การกล่าวหา, คำประนาม, การดำเนิน, การกล่าวโทษ, การใส่ความ
achievement	(อะชีฟ'เวินท์) น. การบรรลุผล, ความสำเร็จ, ความสัมฤทธิ์, ผลสัมฤทธิ์, สัมฤทธิ์ภาพ
acolyte	(แอค'โคไลท) น. เด็กผู้ช่วยฆาตหลวงทำพิธี, พระในนิกายดรมันคาทอลิก
actionon	(แอค'ทอน) น. ธาตุเฉื่อยที่เป็นแก๊สและมีสัญลักษณ์ An; at. no. : 86; at. wt. : 2
actor	(แอค'เทอะ) น. นักแสดงชาย, ผู้กระทำ, ผู้ดำเนินการ (คำที่มีความหมายเหมือนกัน)
acuity	(อะคิว'อิที) น. ความหลักแหลม, ความคม, ความรุนแรง, ความคมกริบ, ความชัดเจน
addend	(แอค'เดนด) น. เลขหรือจำนวนที่บวกเข้าด้วยกัน
adit	(แอค'ดิท) น. ทางตามแนวอนเข้าสู่เหมือง, การเข้าหา, ทางเข้า
admittance	(แอดมิท'เทินซ์) น. การอนุญาตให้เข้า, การรับเข้า, การสารภาพ, การวัดการนำไฟฟ้า
advocacy	(แอค'โวคะซี) น. การเป็นทนาย, ทนาย, การสนับสนุน, ผู้สนับสนุน, การแก้ต่าง
aerodrome	(แเอ'โรโดรม) น. สนามบิน
afflux	(แอฟ'ฟลักซ์) น. สิ่งไหลไปทางจุดหนึ่ง, การไหลไปทาง (that which flows)
afreet	(แอฟ'ริท, อะฟริท') น. ปีศาจร้ายในนิยายอาหรับ (monster, afrit)
ageism	(เอจ'จิสซึม) น. การเลือกที่รักมักที่ชังโดยถืออายุเป็นเกณฑ์และโดยเฉพาะกับคนแก่
agger	(แอก'เจอะ) น. กระแสน้ำที่น้ำลดแล้วค่อยๆ ขึ้นแล้วลดอีก. -S... double tide
agrapha	(แอก'ราฟา) น. คำสอนของพระเยซูคริสต์ที่บันทึกโดยชาวคริสต์ใน New Testa
agrimony	(แอก'กริมอนี) น. ญาพันชุ่มังกรจำพวก Agrimonia
agueweed	(เอ'กิววิด) น. พืชจำพวก Eupatorium ในอเมริกา gentian (Gentiana quinquefolia
aisle	(ไอส์ล) น. ทางเดินระหว่างที่นั่ง (ในรถ, โรงหนัง, เครื่องบิน, โบสถ์). ที่นั่งฝั่งธรรมที่แม่
albania	(แอลเบ'เนีย) น. ประเทศอัลบانيا
album	(แอล'บั้ม) น. สมุดหน้าเปล่าสำหรับเก็บภาพแสตมป์หรืออื่นๆ, แผ่นเสียงขนาดใหญ่, ๑
alcoholicity	(แอลกอฮอล์ลิซ'ซิที) น. คุณภาพหรือความเข้มข้นของเครื่องดื่มแอลกอฮอล์
alexander	(แอลเล็กซาน'เดอะ) น. ชื่อเหล้า Cocktail ชนิดหนึ่ง
alleyway	(แอล'ลีย์เว) น. ทางที่แคบ, ตรอก, ซอย
allium	(แอล'เลียม) น. พืชจำพวก Allium (ต้นหอม, กระเทียม) (onion, leek, shallot)
allograph	(แอล'โลกราฟ) น. การเขียนหรือการเขียนชื่อคนหนึ่งเพื่ออีกคนหนึ่ง

4.2 ส่วนของการทำงานของโปรแกรม

การทำงานจะแบ่งเป็น 2 ส่วน คือส่วนที่เป็น Application ในการค้นหาคำศัพท์และทำ Relation test และส่วนที่เป็น website สำหรับ online dictionary เพื่อค้นหาความหมายของคำศัพท์ และเพื่อค้นหาการทำ relation test

-ในส่วนของ Application เมื่อเริ่มต้น โปรแกรมการทำงานจะปรากฏหน้าจอในรูป 4.1 โดยแบ่งเมนูการทำงานออกเป็นส่วนๆคือส่วนของการเก็บรวบรวม Corpus มาจากเว็บไซต์ ส่วนของการ filtering คำศัพท์มาจาก Corpus ที่เรารวบรวมได้ และส่วนของการทำ Relation test ดังแสดงในรูป 4.1

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.1 หน้าจอหลักของ Application

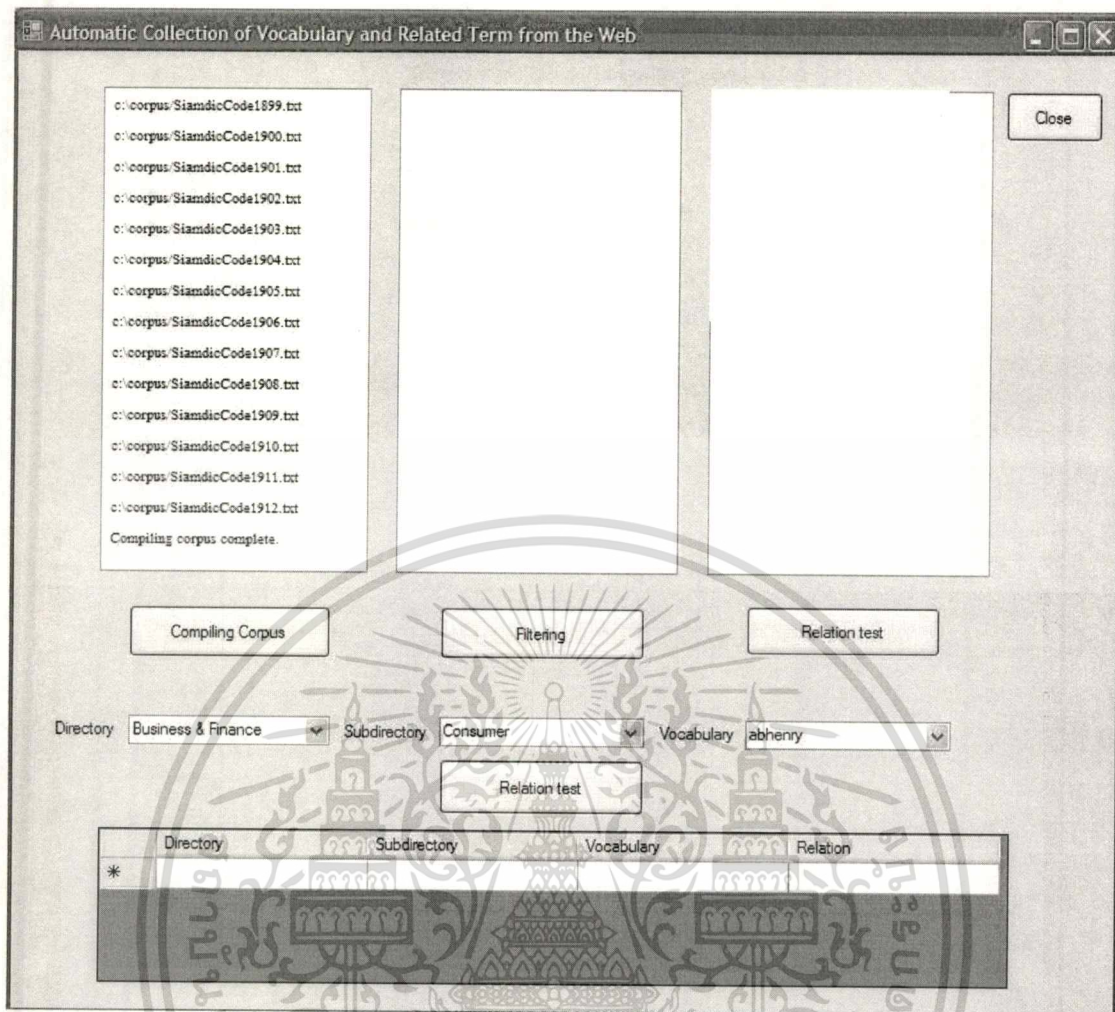
1. ส่วนของการเก็บรวบรวมข้อมูลจาก Corpus จะเป็นการทำงานเมื่อกลุ่ม Compiling Corpus ซึ่งเมื่อทำการกดแล้ว การทำงานคือการเก็บรวบรวมข้อมูลจากเว็บไซต์โดยใช้เทคโนโลยี search engine มาช่วย ซึ่งในส่วนของการทำงานนี้ input ของโปรแกรมส่วนนี้ได้แก่

- URL ของเว็บไซต์ของ dictionary ที่เราคัดเลือกมา
- จำนวนหน้าเว็บไซต์ของเว็บไซต์ dictionary นั้นๆ
- รูปแบบการเก็บ source file ของเว็บไซต์นั้นๆ
- รูปแบบการเก็บคำศัพท์ภาษาอังกฤษของเว็บไซต์นั้นๆ
- รูปแบบการเก็บคำศัพท์ภาษาไทยของเว็บไซต์นั้นๆ

เมื่อทำงานแล้ว ผลของ Corpus ที่ถูกสร้างขึ้นจะโชว์ขึ้นมาให้ดูที่ช่องแสดงผลตามรูป 4.2

โดยที่ถ้าการทำงานเสร็จสมบูรณ์จะขึ้นคำว่า Compiling corpus complete ขึ้นมา ดังรูป 4.2

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



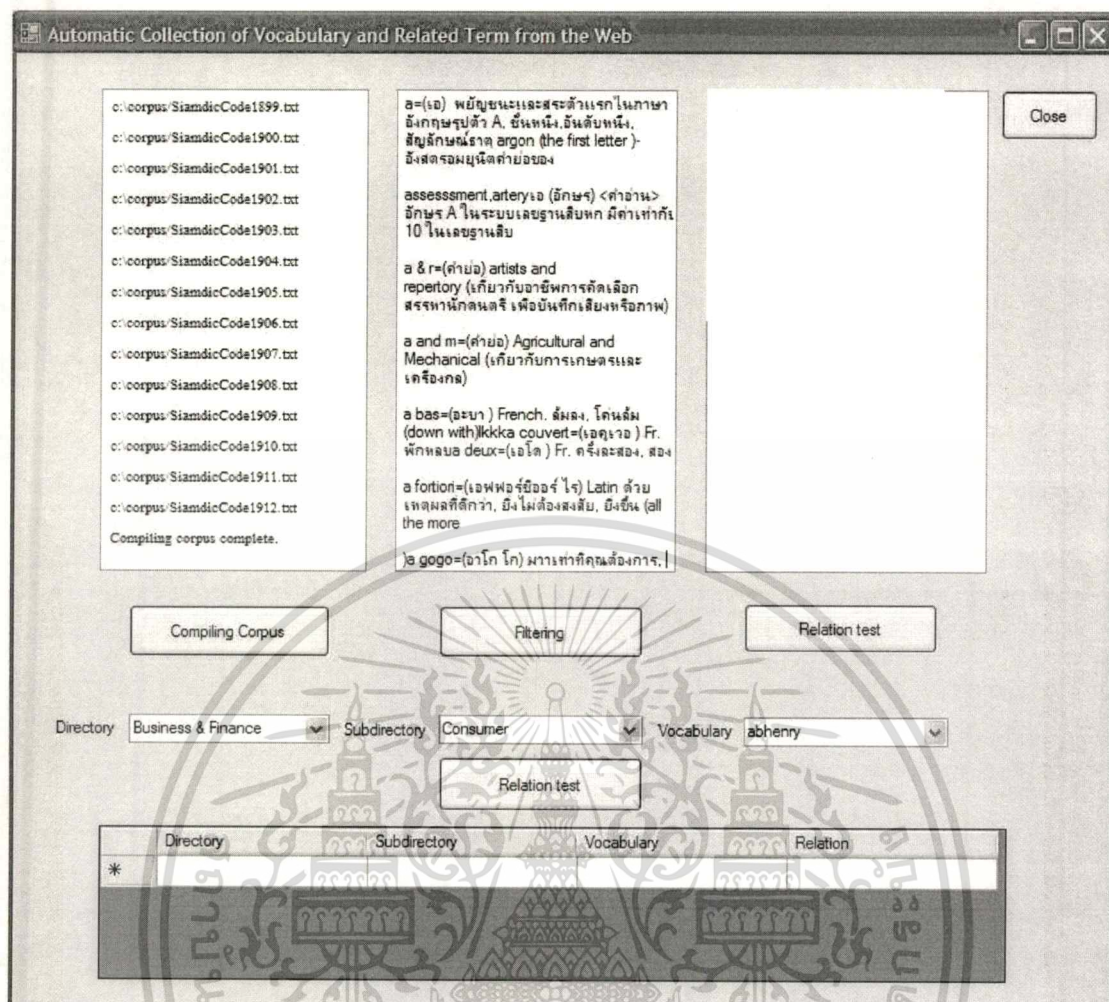
รูปที่ 4.2 ผลลัพธ์จากการกดปุ่ม Compiling corpus ที่ปรากฏบน application

2. ส่วนของการทำการ Filtering เพื่อดึงคำศัพท์และคำแปลออกมาจาก Corpus ที่เรารวบรวมได้จากการทำ compiling corpus โดยการทำงาน 2 ขั้นตอนได้แก่การทำ Web page collection และการทำ Sentence extraction ซึ่งในส่วนของการทำงานนี้ input ของโปรแกรมส่วนนี้ได้แก่

- Corpus ที่เราเก็บรวบรวมได้
- Automatic Term Recognition

เมื่อทำงานแล้ว ผลของ คำศัพท์และคำแปล ที่ถูกทำ Filtering มาได้จะ โฉวขึ้นมาให้ดูที่ช่องแสดงผลตามรูป 4.3 โดยที่ถ้าการทำงานเสร็จสมบูรณ์จะขึ้นคำว่า Filtering complete ขึ้นมา ดังรูป 4.3 และจะทำการเก็บคำศัพท์และความหมายลงในฐานข้อมูล ในตาราง Dictionary

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

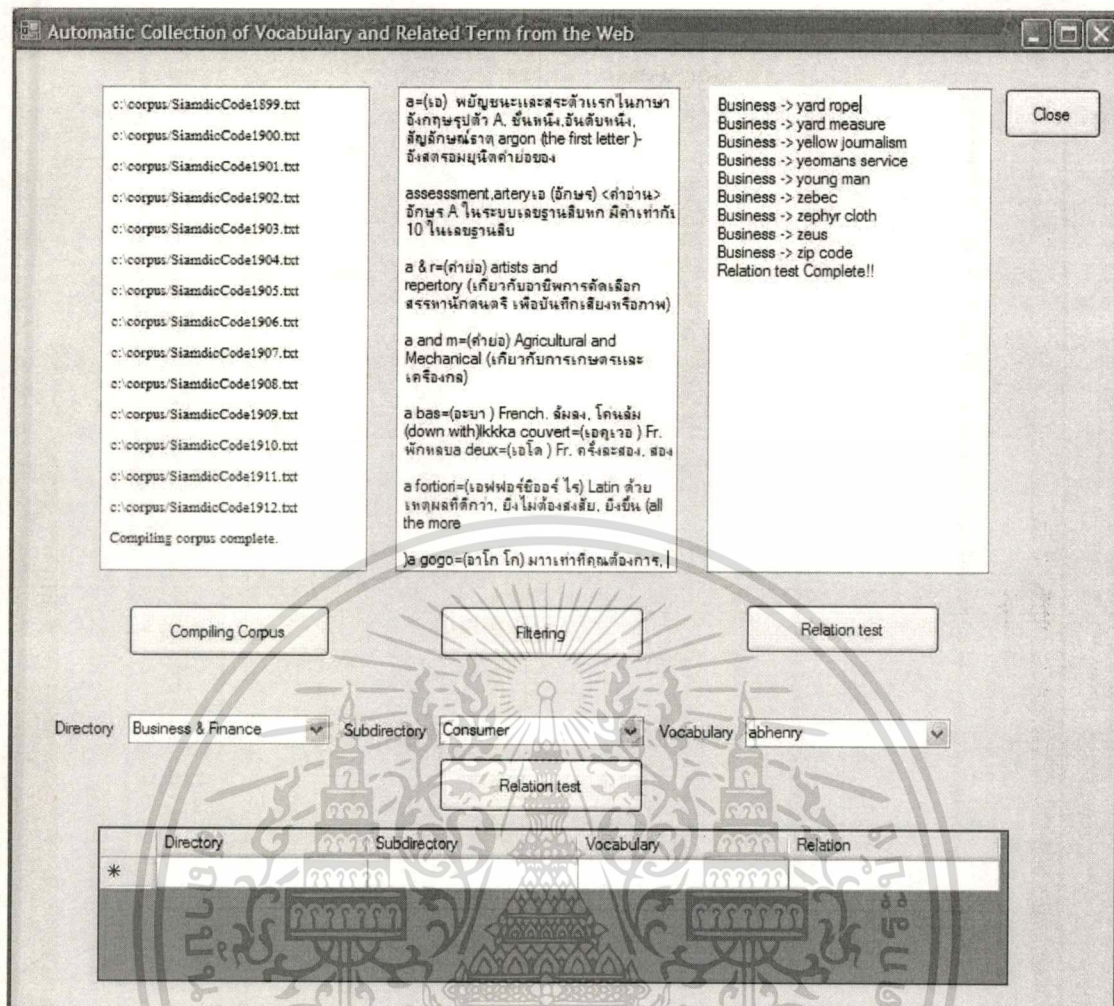


รูปที่ 4.3 ผลลัพธ์จากการกดปุ่ม Filtering ที่ปรากฏบน application

3. ส่วนของการทำการ Relation test เพื่อทำการหาความสัมพันธ์ระหว่างคำศัพท์และหมวดหมู่ โดยจะเริ่มทำงานเมื่อมีการกดปุ่ม Relation test ตัวโปรแกรมจะนำหมวดหมู่ที่เราต้องการ และคำศัพท์ที่เราจัดไว้ มาทำการหาค่า Relation โดยได้ผลเป็นค่าความน่าจะเป็นที่คำศัพท์และหมวดหมู่จะมีความเกี่ยวข้องกัน

เมื่อทำงานแล้ว ผลของความสัมพันธ์ระหว่างคำศัพท์และหมวดหมู่ จากการทำการ Relation test จะแสดงขึ้นมาให้ดูที่ช่องแสดงผลตามรูป 4.4 โดยที่ถ้าการทำงานเสร็จสมบูรณ์จะขึ้นคำว่า Relation complete ขึ้นมา ดังรูป 4.4

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.4 ผลลัพธ์จากการกดปุ่ม Relation ที่ปรากฏบน application

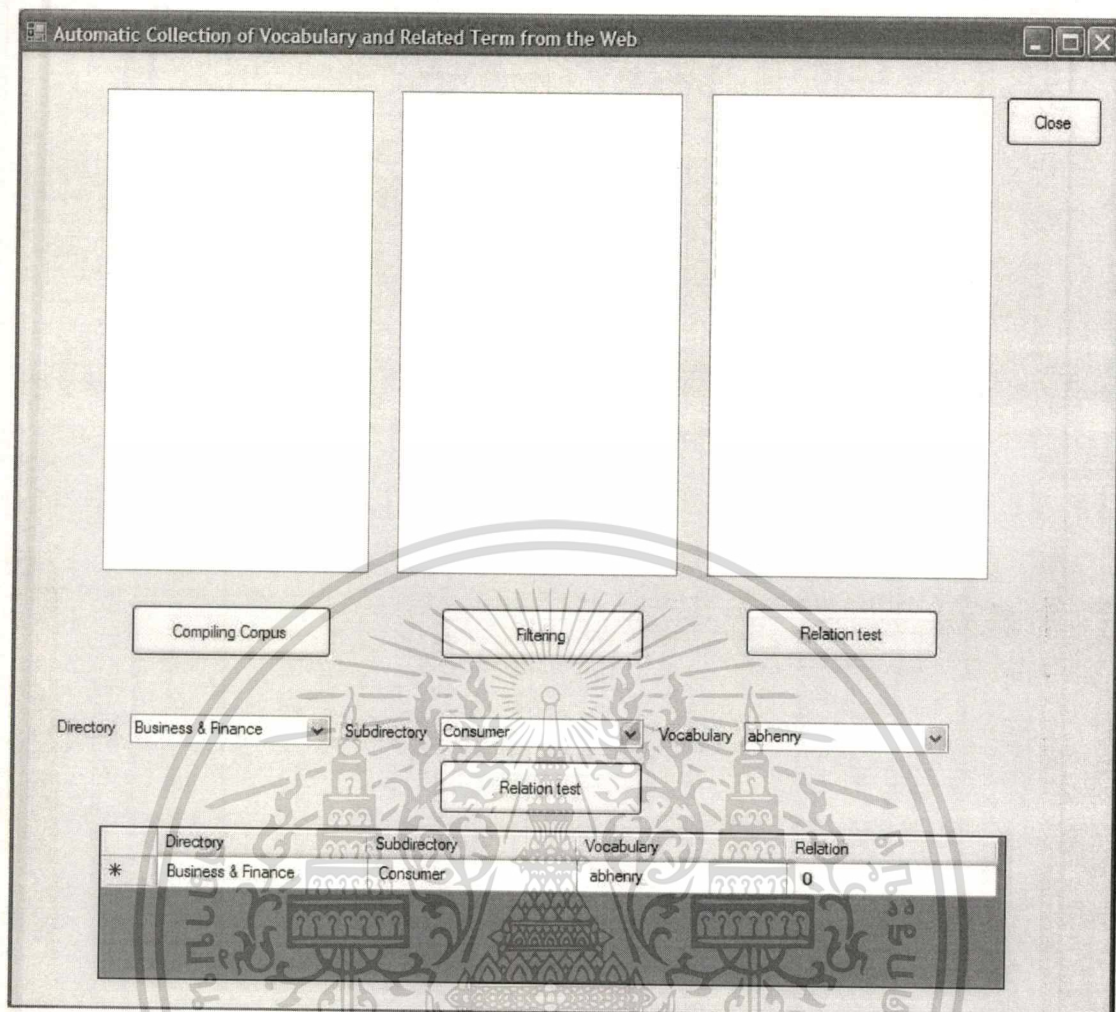
และผลของการทำ Relation test จะ ได้ผลดังตาราง 4.4

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.4 แสดงผลที่ได้จากการทดลองของ Relation test

relation_id	Word_id	Subdirectory_id	Directory_id	probability	Category_hit	Word_hit	Both_hit
210301	12541	9	2	0	136000000	18	0
210302	12541	11	2	0	1150000000	18	0
210303	12541	12	2	0.1111111111111111	474000000	18	2
210304	12541	14	2	0.0555555555555556	6450000	18	1
210305	12541	16	3	0.0555555555555556	975000000	18	1
210306	12541	19	3	0	209000000	18	0
210307	12541	21	3	0.3333333333333333	149000000	18	6
210308	12541	23	3	0	1150000	18	0
210309	12541	27	3	0.0555555555555556	35000000	18	1
210310	12541	28	3	0.2222222222222222	863000000	18	4
210311	12541	34	3	0.1111111111111111	1010000000	18	2
210312	12541	36	1	0.4444444444444444	223000000	18	8
210313	12541	37	1	0.2222222222222222	530000000	18	4
210314	12541	38	1	0.0555555555555556	287000000	18	1
210315	12541	39	1	0	265000000	18	0
210316	12541	40	1	0.0555555555555556	83800000	18	1
210317	12541	42	1	0.4444444444444444	2460000000	18	8
210318	12541	43	1	0.1666666666666667	60500000	18	3
210319	12541	46	1	0.1111111111111111	248000000	18	2
210320	12541	53	4	0	1380000	18	0
210321	12541	55	4	0.1111111111111111	9010000	18	2
210322	12541	57	4	0	5340000	18	0
210323	12541	60	4	0	295000000	18	0
210324	12541	66	4	0.0555555555555556	516000000	18	1

4. ส่วนของการเลือกทดสอบ Relation test แบบเจาะจงคำศัพท์หรือตามกลุ่มของตัวอักษร เพื่อหาความสัมพันธ์ระหว่างหมวดหมู่ที่เลือก หรือทุกหมวดหมู่ กับคำศัพท์ที่เรากำหนดไว้ โดยหลักของการทำงานนี้ ใช้เพื่อตรวจสอบความสัมพันธ์ของคำศัพท์เฉพาะเจาะจงที่ต้องการ โดยการเลือกหมวดหมู่ที่ต้องการ เลือกหมวดหมู่ย่อยที่ต้องการ เลือกคำศัพท์ที่ต้องการ หรือเลือกคำขึ้นต้นที่ต้องการ เมื่อเลือกทั้งหมดเสร็จสิ้น การทำงานจะเริ่มเมื่อทำการกดปุ่ม Relation test ผลลัพธ์ของการทำ Relation test จะแสดงให้เห็นในช่องแสดงผลบน application ดังรูป 4.5



รูปที่ 4.5 แสดงผลที่ได้จากการทำ Relation test แบบเจาะจงคำศัพท์

-ในส่วนของ online dictionary เราจะแสดงผลของการทำงานของโปรแกรมตามทฤษฎี Automatic Collection of Vocabulary and Related Term from the Web ด้วยการ ใช้ website ในที่นี้ เราจัดทำในรูปแบบของ online dictionary โดยใช้ชื่อว่า ART Dictionary ซึ่งจะมีการทำงาน 3 รูปแบบ ได้แก่

1. การทำงาน Online Dictionary ทั่วไป
2. การทำงานของการตรวจสอบ Relation test ตามหมวดหมู่
3. การเลือกคำศัพท์ที่ต้องการเพื่อตรวจสอบผลของ Relation test กับทุกๆหมวดหมู่ย่อย

-การทำงานของ Dictionary ทั่วไปโดยหน้าแรกของ online dictionary สามารถดูได้ดังรูป 4.6 ซึ่งเป็นความสามารถในการค้นหาความหมายของคำศัพท์ในรูปแบบทั่วไป โดยแบ่งการทำงานในส่วนนี้ออกเป็น 3 ส่วนย่อยได้แก่
เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์การใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

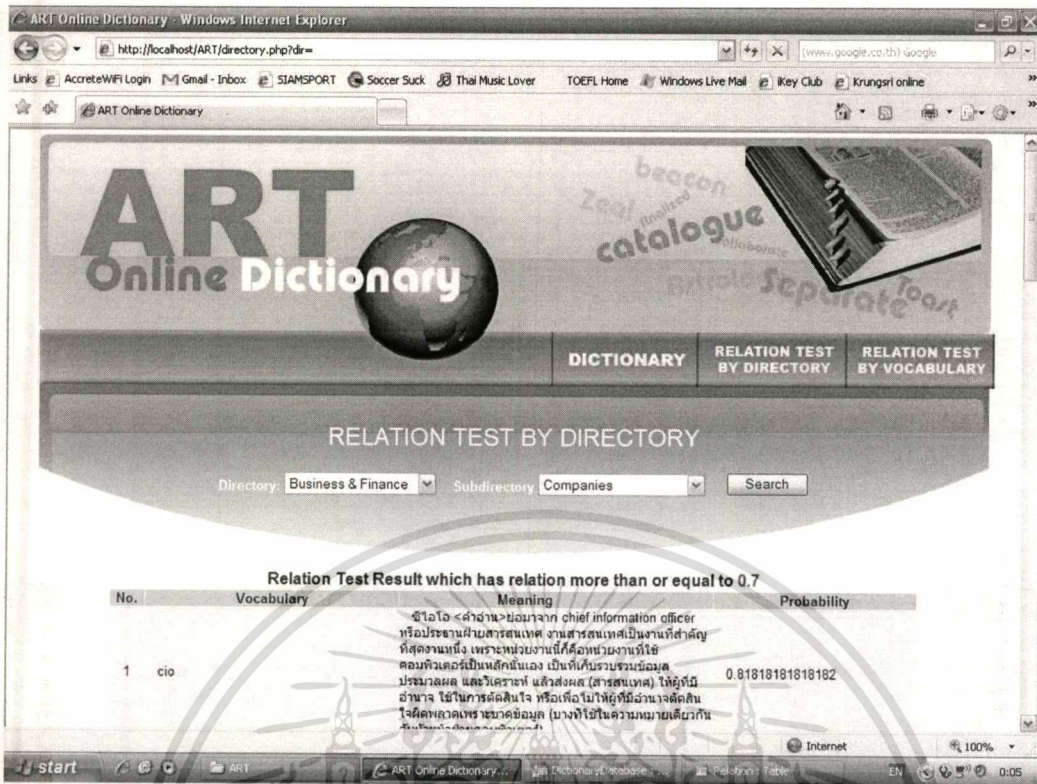
1. การค้นหาคำศัพท์เฉพาะคำศัพท์ที่เราต้องการ
2. การค้นหาคำศัพท์ที่มีคำศัพท์ที่เราต้องการเป็นตัวประกอบ
3. การค้นหาคำศัพท์ที่ขึ้นต้นด้วยคำศัพท์ที่เราต้องการ



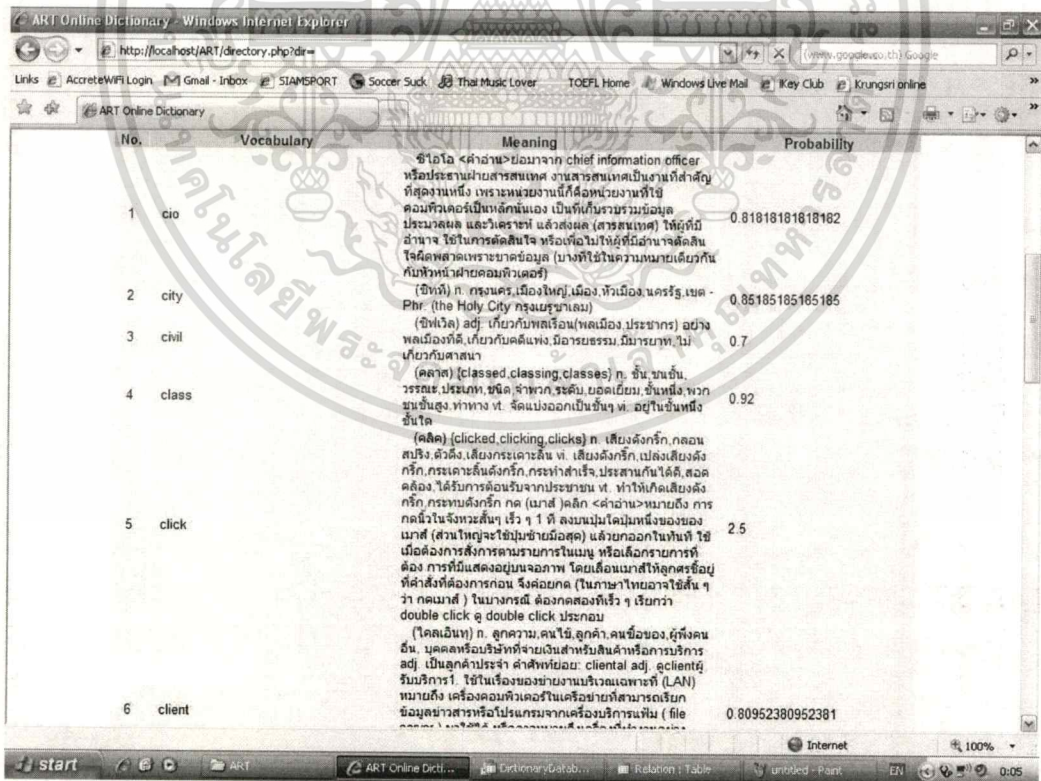
รูปที่ 4.6 แสดงรูปแบบหน้าการทำงานหลักของ online dictionary

-การทำงานของการตรวจสอบ Relation test ตามหมวดหมู่ จะเป็นการเลือกดูค่า Relation test โดยเลือกจากหมวดหมู่หลัก และหมวดหมู่ย่อย และ online dictionary จะแสดงผลของคำศัพท์ต่างๆที่มีค่า Relation test กับหมวดหมู่ย่อยที่ผู้ใช้งานเลือก โดยที่ผลลัพธ์ที่จะมีจำนวนเท่าใด หรือครอบคลุมการคำศัพท์ทุกคำหรือไม่ ขึ้นกับการทำงานของส่วน application ซึ่งสามารถจะหาค่า Relation test ของคำศัพท์ทุกคำ หรือแค่บางส่วนได้ โดยผลของการทำงานสามารถดูได้ดังรูป 4.7 และรูป 4.8 ซึ่งเป็นความสามารถในการตรวจสอบ Relation test ตามหมวดหมู่

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีกรุณาไปใช้



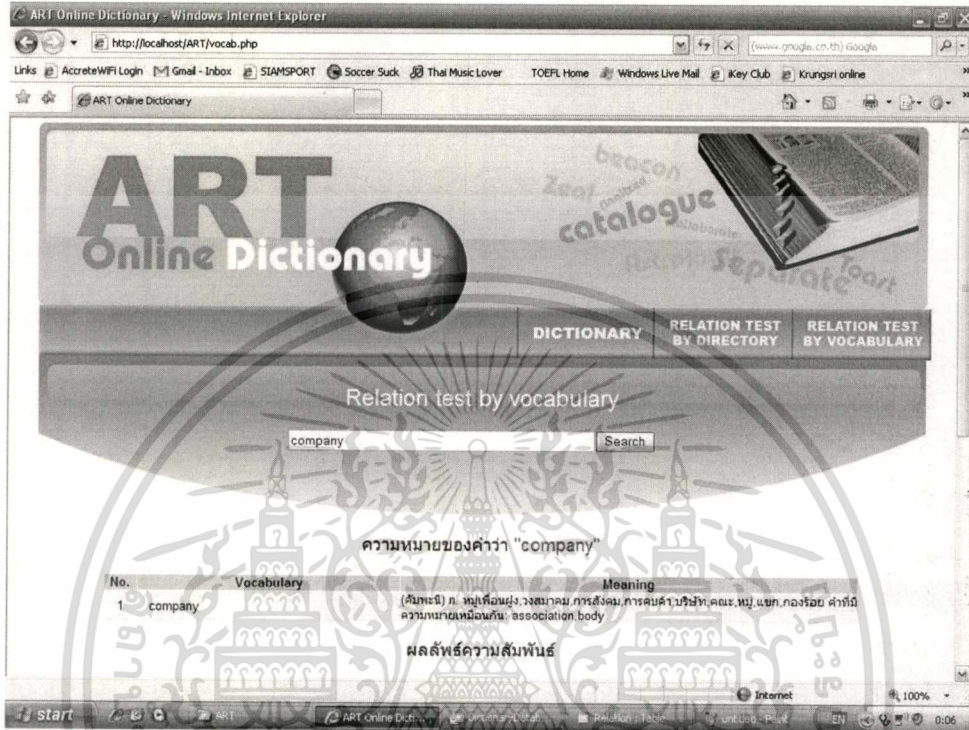
รูปที่4.7 แสดงการทำงานของ การตรวจสอบ Relation test ตามหมวดหมู่



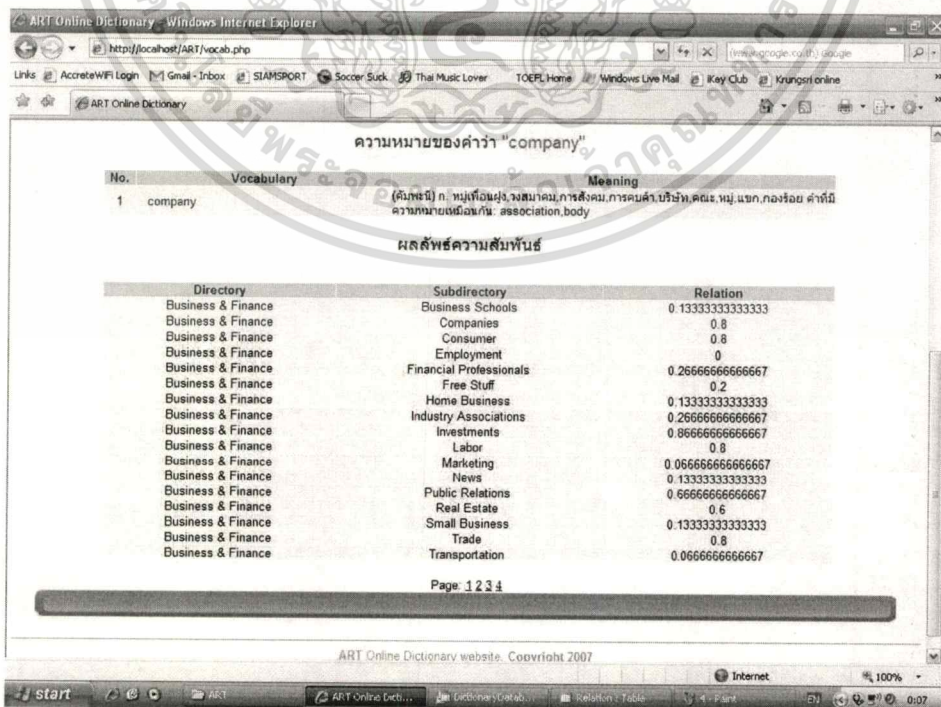
รูปที่4.8 แสดงการทำงานของ การตรวจสอบ Relation test ตามหมวดหมู่ (เพิ่มเติม)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- การเลือกคำศัพท์ที่ต้องการเพื่อตรวจสอบผลของ Relation test กับทุกๆหมวดหมู่ย่อย โดยผลของการทำงานสามารถดูได้ดังรูป 4.9 และรูป 4.10 ซึ่งเป็นความสามารถในการตรวจสอบ Relation test ของทุกๆหมวดหมู่ตามคำศัพท์ที่เราต้องการ โดยที่ผลลัพธ์ที่ได้จะมีการแสดงความหมายของคำศัพท์นั้นๆก่อน และมีการแสดงค่า Relation test กับหมวดหมู่ต่างๆต่อเนื่องลงมา



รูปที่ 4.9 แสดงการทำงานของ การตรวจสอบ Relation test ตามคำศัพท์



รูปที่ 4.10 แสดงการทำงานของ การตรวจสอบ Relation test ตามคำศัพท์ (เพิ่มเติม)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่นับผูกมัดให้เข้าไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีก้ารนำไปใช้

จากผลการทดลองที่ได้ทำใน โปรแกรมนี้ เราจะแบ่งผลการทดลองออกเป็นส่วนๆ ตามค่าของ Relation test ที่ได้เป็น 4 ชนิดดังนี้

-ชนิดที่ 1 ค่า Probability =1 จะแบ่งได้ 2 ประเภท คือ

Type 0: หมวดหมู่และคำศัพท์อาจจะเป็นคำๆเดียวกัน

Type 1: คำศัพท์ที่มีความเกี่ยวข้องกับหมวดหมู่หมวดหมู่มาก ดังตัวอย่างรูป 4.11

ความหมายของคำว่า "Company"

Subdirectory	Relation
Labor	1.0
Trade	1.0

รูปที่ 4.11 แสดงผลของ Relation test สำหรับคำศัพท์ที่มีความเกี่ยวข้องกับหมวดหมู่หมวดหมู่

-ชนิดที่ 2 ค่า Probability <1 และค่า Probability ≥ 0.8

Type 2: คำศัพท์ที่มีความเกี่ยวข้องกับหมวดหมู่ ดังตัวอย่างรูป 4.12

ความหมายของคำว่า "Company"

Subdirectory	Relation
Companies	0.8
Consumer	0.8

รูปที่ 4.12 แสดงผลของ Relation test สำหรับ คำศัพท์ที่มีความเกี่ยวข้องกับหมวดหมู่

-ชนิดที่ 3 ค่า Probability <0.8 และค่า Probability ≥ 0.5

Type 3: คำศัพท์และหมวดหมู่อาจเกี่ยวข้องหรือไม่เกี่ยวข้องกัน ดังตัวอย่างรูป 4.13

ความหมายของคำว่า "connection"

Subdirectory	Relation
Companies	0.5
Consumer	0.75

รูปที่ 4.13 แสดงผลของ Relation test สำหรับคำศัพท์และหมวดหมู่อาจเกี่ยวข้องหรือไม่เกี่ยวข้องกัน

-ชนิดที่ 4 ค่า Probability < 0.5

Type 4: คำศัพท์และหมวดหมู่ไม่มีความเกี่ยวข้องกันเลย ดังตัวอย่างรูป 4.14

ความหมายของคำว่า "connection"

Subdirectory	Relation
Financial Professionals	0.041666666666667
Free Stuff	0.166666666666667

รูปที่ 4.14 แสดงผลของ Relation test สำหรับคำศัพท์และหมวดหมู่ไม่มีความเกี่ยวข้องกันเลย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 5

สรุปผลและข้อเสนอแนะ

จากการศึกษาการทำงานของการทำงานการเก็บรวบรวมคำศัพท์และคำแปลจากหน้าเว็บไซต์และการค้นหาความสัมพันธ์ของคำศัพท์และหมวดหมู่ต่าง โดยใช้ทฤษฎี Automatic Collection of Vocabulary and Related Term from the Web ในโครงการนี้ สามารถสรุปผลการดำเนินงานและสรุปผลการทดลอง รวมถึงข้อเสนอแนะได้ดังต่อไปนี้

5.1 สรุปผลการทำงาน

จากการศึกษาใน โครงการนี้สามารถสรุปผลการดำเนินงานในการทำ Automatic Collection of Vocabulary and Related Term from the Web ได้ดังนี้

- ในการเก็บรวบรวมข้อมูลคำศัพท์และคำแปลมาจากเว็บไซต์ สามารถทำได้ 2 รูปแบบ คือ การเก็บรวบรวมคำศัพท์จากเว็บไซต์ทั่วไป โดยใช้เทคโนโลยี search engine โดยที่การเก็บรวบรวมคำศัพท์จากเว็บไซต์ทั่วไป จะได้คำศัพท์ประมาณ 3 คำจากหน้าเว็บไซต์ประมาณ 3000 เว็บไซต์ ซึ่งเป็นจำนวนน้อยมากและใช้เวลานานในการรวบรวมคำศัพท์ สำหรับกรณีที่ต้องการสามารถให้ program ทำงานโดยใช้เวลานานได้ แต่ในกรณีที่ต้องการคำศัพท์จำนวนมากในเวลาจำกัด อีกทางเลือกคือการค้นหาคำศัพท์จาก website online dictionary ที่มีอยู่ทั่วไป เนื่องจากในเว็บไซต์เหล่านี้มีคำศัพท์และคำแปลจำนวนมาก ทำให้เราสามารถรวบรวมข้อมูลได้ แต่ในกรณีที่มีเวลาในการทำงาน ควรจะทำการรวบรวมคำศัพท์จากเว็บไซต์ทั่วไปเพิ่มเติมจากคำศัพท์ที่ได้จาก online dictionary ปรกติเพื่อการได้คำศัพท์ที่มีความหลากหลายมากขึ้น

- ในการทำ Compiling Corpus ใช้เวลานานในการรวบรวมคำศัพท์ เนื่องจากมีคำศัพท์ในเว็บไซต์จำนวนมาก ทำให้เครื่องคอมพิวเตอร์ที่จะใช้ในการทำงานนี้ต้องมีประสิทธิภาพสูงและสามารถทำงานต่อเนื่องได้เป็นเวลานาน จึงจะสามารถทำการรวบรวมคำศัพท์ได้ โดยควรแบ่งการรวบรวมคำศัพท์ออกเป็นส่วนๆ เช่นการแบ่งตามตัวอักษร เพื่อให้เครื่องคอมพิวเตอร์ไม่ต้องทำงานมากจนเกิดความสามารถและแบ่งการทำงานได้

- ในการทำการ Filtering คำศัพท์ ด้วยเทคนิค Web page collection และการทำ Sentence extraction ใช้เวลานานในการคัดแยกคำศัพท์ที่ออกมาจาก Corpus ที่รวบรวมไว้ ซึ่งใช้เวลามากกว่า การทำ Compiling Corpus ดังนั้นก็เช่นกับการทำ Compiling Corpus คือควรมีการแบ่งตามตัวอักษรเพื่อการทำงานที่มีประสิทธิภาพ หรือต้องใช้เครื่องคอมพิวเตอร์ที่มีประสิทธิภาพเพียงพอในการทำงาน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- ในการทำ Relation test เป็นส่วนของการทำงานที่ใช้เทคโนโลยี search engine เข้าช่วย ดังนั้นเครื่องคอมพิวเตอร์ที่ใช้ในการคำนวณต้องมีการต่ออินเทอร์เน็ตตลอดเวลา และด้วยการที่เรามีคำศัพท์และหมวดหมู่ย่อยจำนวนมากทำให้เราจะเป็นต้องใช้อินเทอร์เน็ตความเร็วสูง

- ผลที่ได้จากการทำ Relation test ยังคงต้องใช้การทำ manual checking ในการดูว่าผลที่ได้ นั้นควรจะสรุปผลอย่างไรซึ่งผลได้ที่ปริมาณมากทำให้ต้องใช้เวลาจำนวนมากในการวิเคราะห์ผล

- เมื่อมีการเรียนรู้จากผลของการทำ Relation test ทำให้ต้องมีการปรับปรุงการทำงาน Relation test ให้มีประสิทธิภาพยิ่งขึ้น โดยใช้เทคนิคต่างๆของ Automatic Collection of Vocabulary and Related Term มาปรับปรุงเช่นการคัดเลือกเฉพาะคำศัพท์ที่เป็นชนิดคำนามในการหา Relation test เนื่องจากจะช่วยเพิ่มความแม่นยำของผลลัพธ์ และมีการปรับเปลี่ยนสมมุติฐานที่เราจะได้ในตอนแรกและข้อสรุปที่ได้หลังจากได้ผลลัพธ์ออกมาแล้ว

5.2 สรุปผลการทดลอง

1. คำที่ได้จากการทำ Relation test ขึ้นกับคำนาม ชนิดของคำนาม และหมวดหมู่ที่จัดเตรียมไว้ คำนามนั้นจะมีความเกี่ยวข้องกับหมวดหมู่น้อยแค่ไหน หรือตรงตามที่ต้องการหรือไม่ขึ้นอยู่กับการจัดเตรียมคำศัพท์และการหมวดหมู่ต่าง ๆ นั้นเอง

2. คำ Relation test ที่ได้จะมีประโยชน์ต่อการหาคำศัพท์ที่มีความเกี่ยวข้องกัน หรือคำศัพท์กับหมวดหมู่ที่มีความเกี่ยวข้องกัน เพื่อใช้ประโยชน์ต่างๆมากมาย เช่นการทำวิจัย การเขียน วิจารณ์งานภาษาอังกฤษ หรือการนำไปทำ machine translation ต่อไปได้

3. การหาคำ Relation test เพื่อให้ได้ผลที่ถูกต้องแม่นยำยิ่งขึ้น ต้องมีการพัฒนาปรับปรุงวิธีการทำ Relation test และต้องมีการคัดเลือกคำศัพท์และหมวดหมู่ให้ตรงกับความต้องการในการใช้งานให้มากที่สุด และผลลัพธ์ที่ได้จะตรงตามความต้องการของผู้ใช้งานมากยิ่งขึ้น

5.3 ข้อเสนอแนะ

เพื่อปรับปรุงโปรแกรมประยุกต์นี้ในอนาคต ผู้ศึกษามีความเห็นว่ เพื่อเพิ่มความเร็วให้กับการทำงานของ Automatic Collection of Vocabulary and Related Term from the Web คือการเตรียม Hardware และสถานะแวดล้อมที่เหมาะสมกับการทำงาน และถ้ามีเวลานานสำหรับการทำ Relation test ควรจะทำ Relation test กับคำศัพท์ทุกตัว เพื่อนำมาใช้ประโยชน์อื่นๆและเพื่อวิเคราะห์และปรับปรุงการทำ Relation test ให้ได้ผลลัพธ์ที่ดียิ่งขึ้น

โปรแกรมประยุกต์นี้อาจจะนำไปใช้พัฒนาต่อเนื่องเพื่อการทำ Relation test ที่มีประสิทธิภาพมากขึ้น หรือการนำไปใช้ในการทำงานอื่นๆเช่น machine translation เป็นต้น

บรรณานุกรม

- ไพศาล โมลิสกุลมงคล. 2545. **Microsoft Visual C#.net**. กรุงเทพฯ : หจก. ไทยเจริญการพิมพ์.
 ศุภชัย สมพานิช. 2546. **คู่มือการเขียนโปรแกรม Visual C#.NET ฉบับโปรแกรมเมอร์**. นนทบุรี
 อินโฟเพรส.
 สิริศักดิ์ คล่องดี. 2542. **สร้างฐานข้อมูลด้วย Microsoft Access2000 อย่างมืออาชีพ**. กรุงเทพฯ
 : เบลโล่การพิมพ์(1988).
 Sato and Y.Sasaki. 2003. "Automatic collection of related terms from the web." 121-124. In
 Proc



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



ภาคผนวก

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



ภาคผนวก ข.

คู่มือการใช้งานระบบ

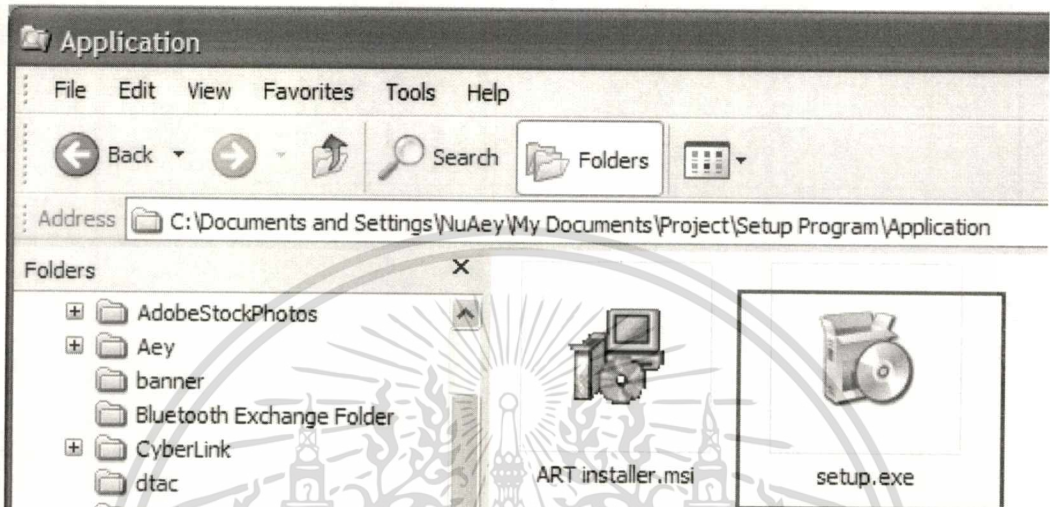
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

คู่มือการติดตั้งระบบการหาค่าความสัมพันธ์ระหว่างคำศัพท์กับหมวดหมู่

1. การติดตั้งโปรแกรมส่วนวินโดวส์แอปพลิเคชัน

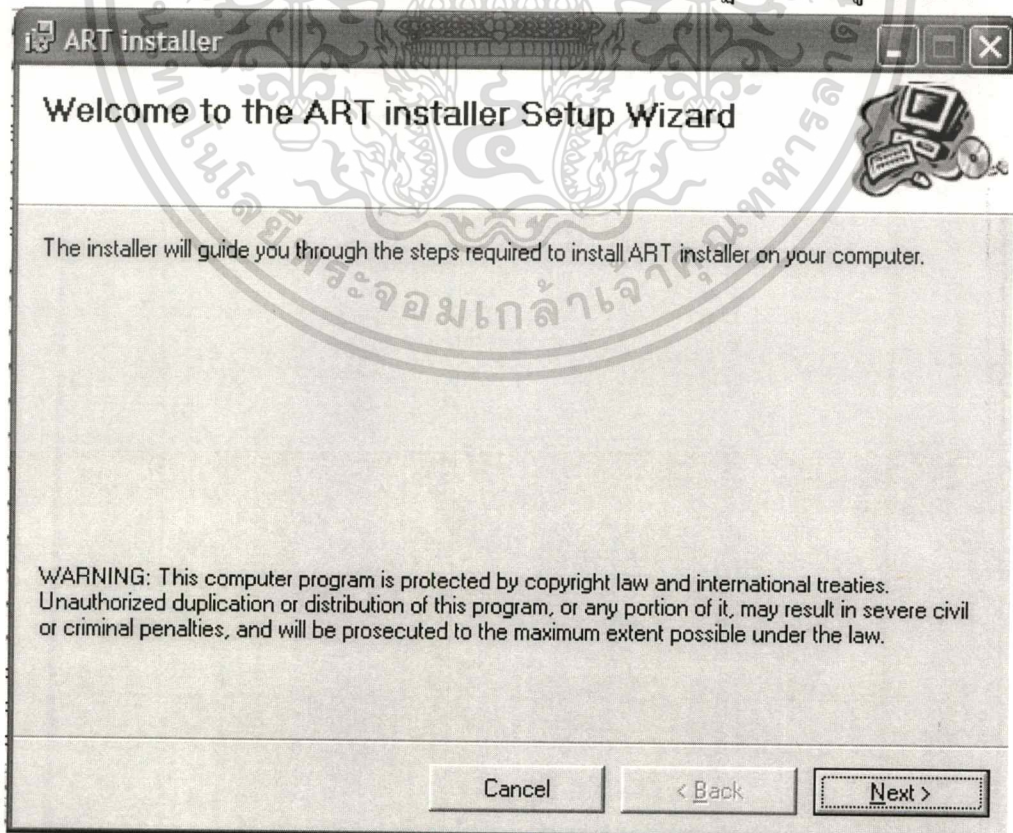
โดยมีขั้นตอนเป็นไปตาม Step by Step ดังต่อไปนี้

1.1 ดับเบิลคลิกที่ตัวไอคอน setup.exe ตามรูปที่ ก.1 ต่อไปนี้



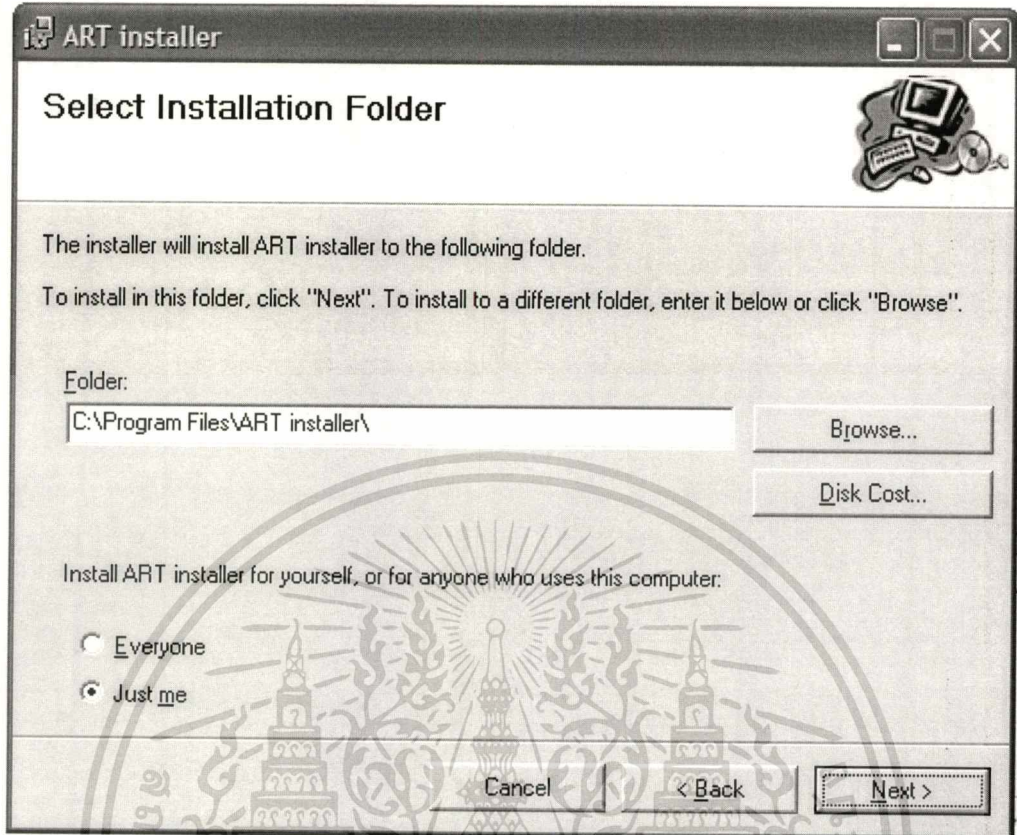
รูปที่ ก.1 หน้าจอแสดงไอคอนตัวติดตั้งโปรแกรม

1.2 เมื่อดับเบิลคลิกที่ไอคอนตัวติดตั้งเรียบร้อยแล้ว จะปรากฏหน้าจอ ดังรูปที่ ก.2 นี้



เอกสารนี้เป็นเอกสารที่สงวนไว้รูปที่ ก.2 หน้าจอแสดงขั้นตอนการติดตั้งโปรแกรมนำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1.3 หลังจากนั้นให้คลิก Next จะปรากฏหน้าจอดังรูปที่ ก.3 นี้



รูปที่ ก.3 หน้าจอแสดงการกำหนดพาธการติดตั้งโปรแกรม

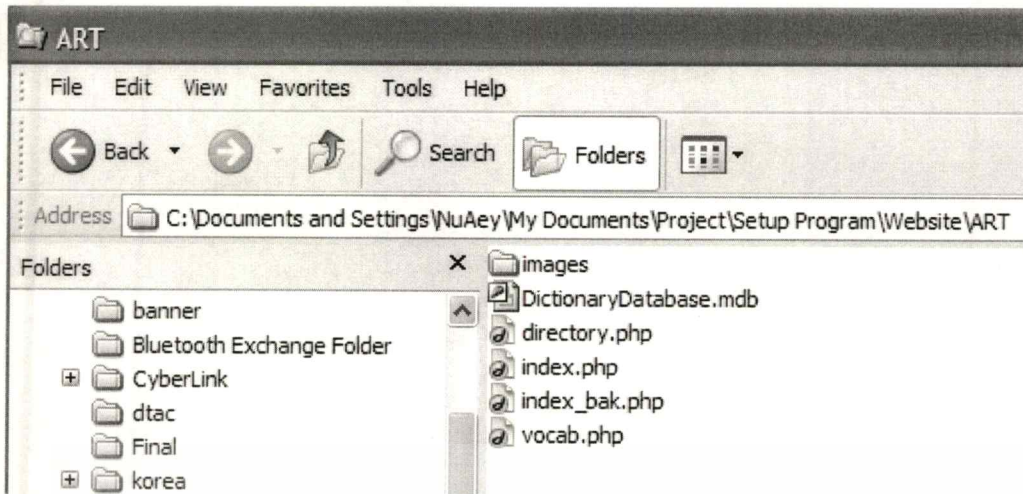
1.4 หลังจากกำหนดพาธการติดตั้งโปรแกรมเรียบร้อยแล้ว คลิก Next ต่อไปจนเสร็จสิ้นขั้นตอนการติดตั้งโปรแกรม จะปรากฏไอคอนโปรแกรมบนหน้าจอเดสทอป ดังรูปที่ ก.4 นี้



รูปที่ ก.4 หน้าจอแสดงไอคอนของโปรแกรม

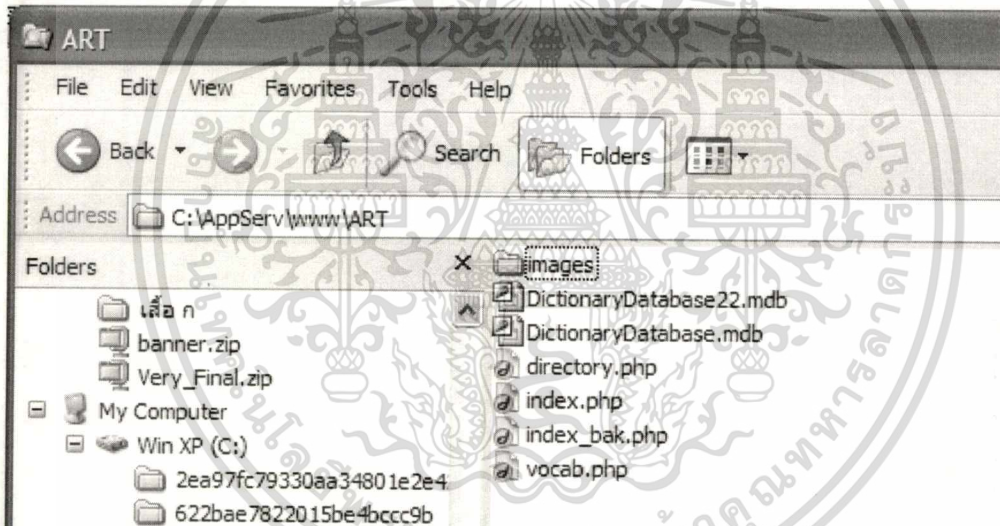
2. การติดตั้งโปรแกรมส่วนเว็บแอปพลิเคชัน

2.1 นำ file.php ทั้งหมดใน folder website ตามรูปที่ ก.5 ต่อไปนี้



รูปที่ ก.5 หน้าจอแสดง file.php ทั้งหมดที่นำมาใช้งาน

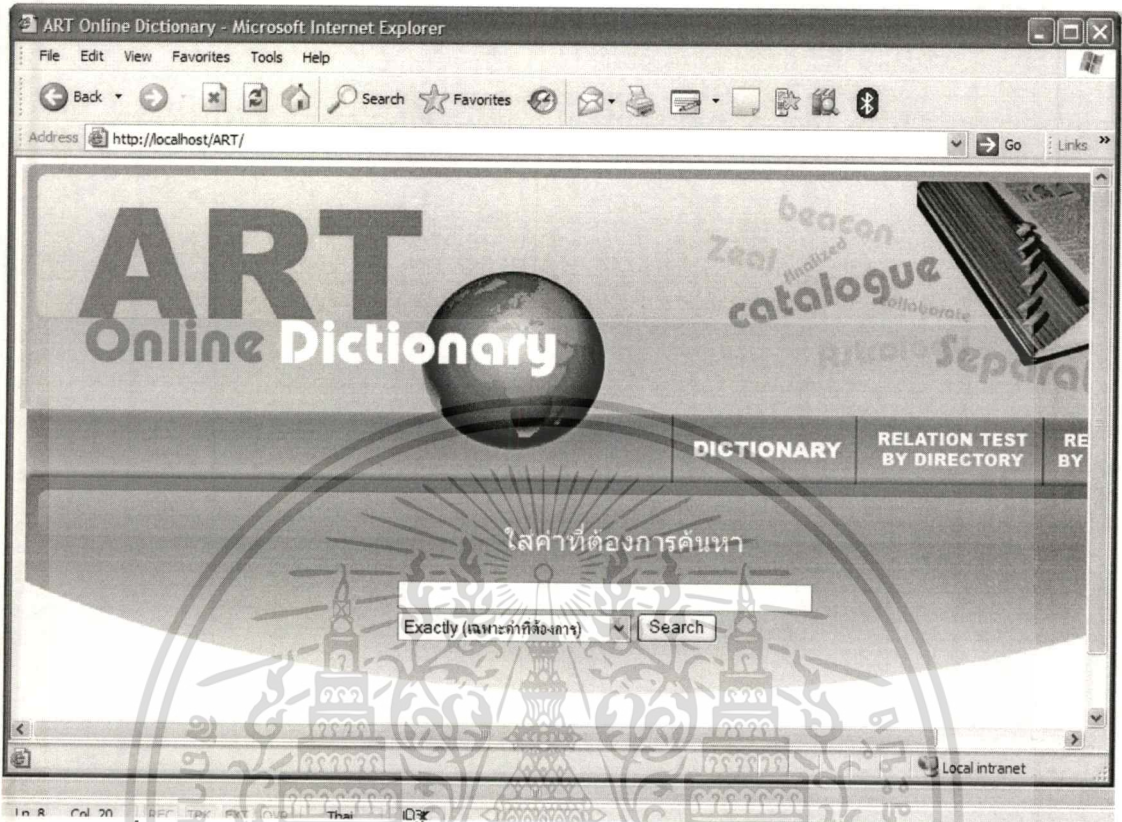
2.2 ติดตั้งโปรแกรม Appserv ให้เสร็จสมบูรณ์ และนำ file ทั้งหมดวางใน Folder www ซึ่งอยู่ภายใต้ Folder AppServ ตามรูปที่ ก.6 ต่อไปนี้



รูปที่ ก.6 หน้าจอแสดง file ทั้งหมดที่นำมาวางใน path ที่เกิดจากโปรแกรม Appserv

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.3 เรียกดู website ด้วย URL “http://localhost/ART ตามรูปที่ ก.7 ต่อไปนี้



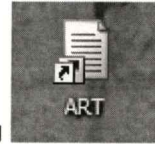
รูปที่ ก.7 หน้าจอแสดง file ทั้งหมดที่นำมาวางใน path ที่เกิดจากโปรแกรม Appserv

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

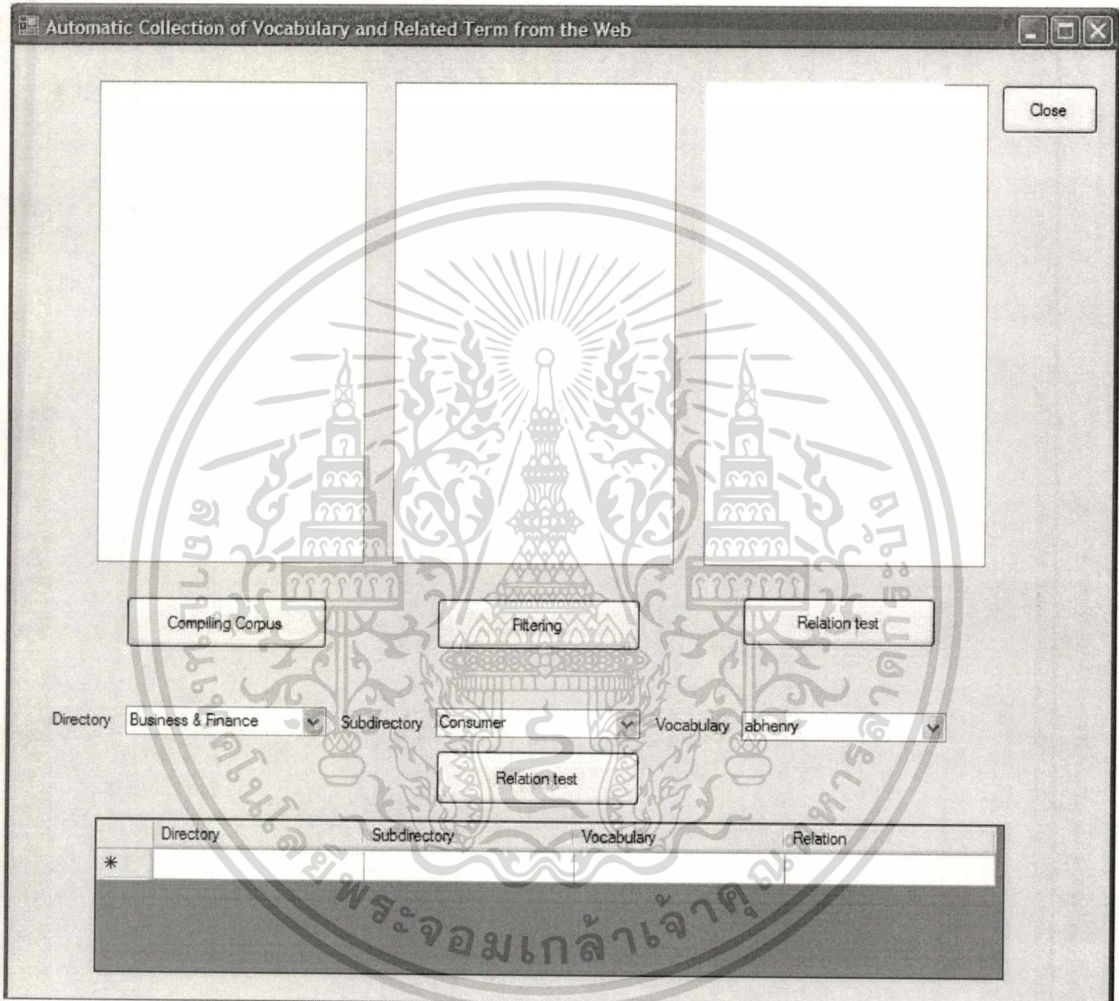


เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

คู่มือการใช้งานระบบการหาค่าความสัมพันธ์ระหว่างคำศัพท์กับหมวดหมู่



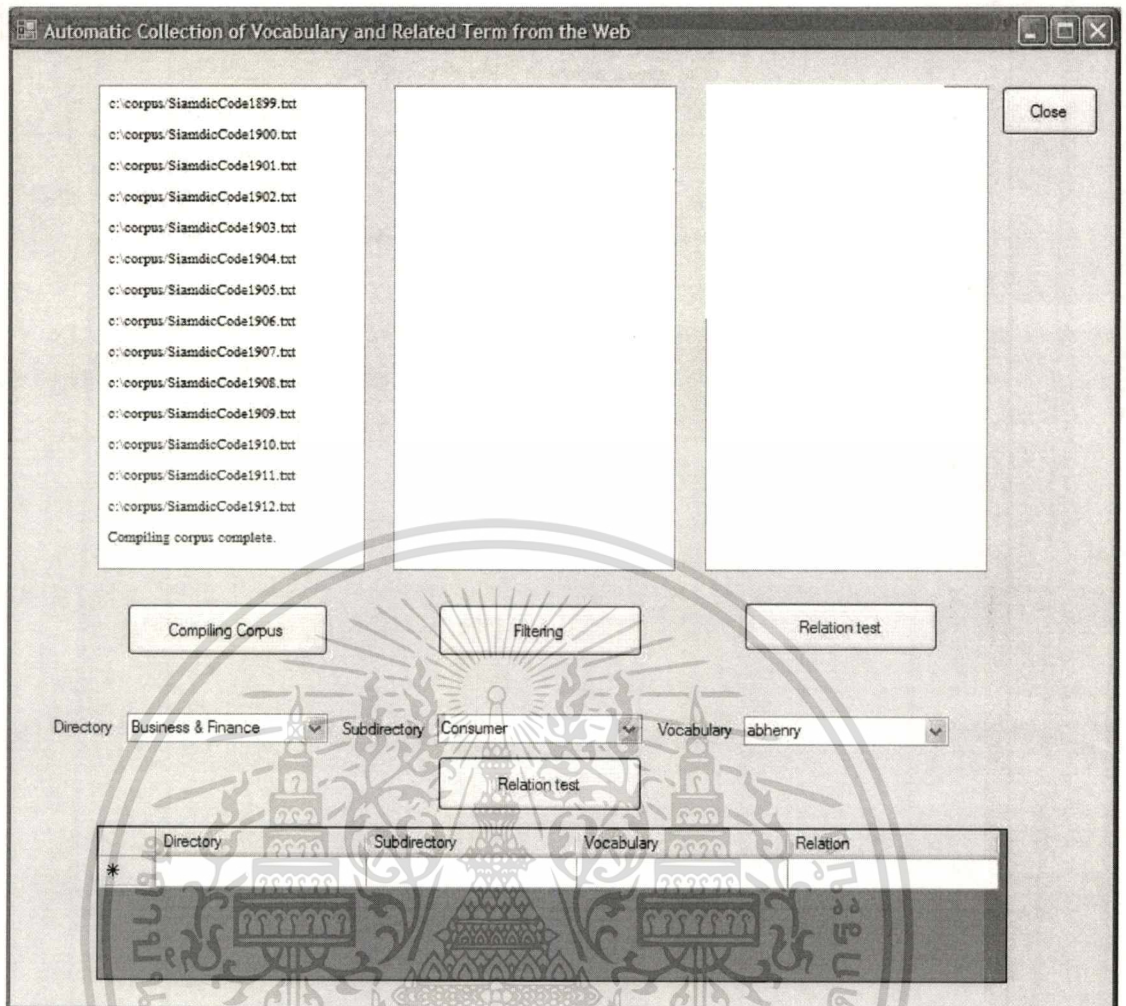
1. เข้าสู่โปรแกรมด้วยการดับเบิลคลิกที่ไอคอนโปรแกรม จะปรากฏหน้าจอเมนูหลักดังรูป ข.8 นี้



รูปที่ ข.8 หน้าจอเมนูหลัก

2. กดปุ่ม Compiling Corpus ซึ่งเมื่อทำการกดแล้ว การทำงานคือการเก็บรวบรวมข้อมูลจากเว็บไซต์ เมื่อการทำงานเสร็จสมบูรณ์จะปรากฏหน้าจอเมนูหลักดังรูป ข.9 นี้

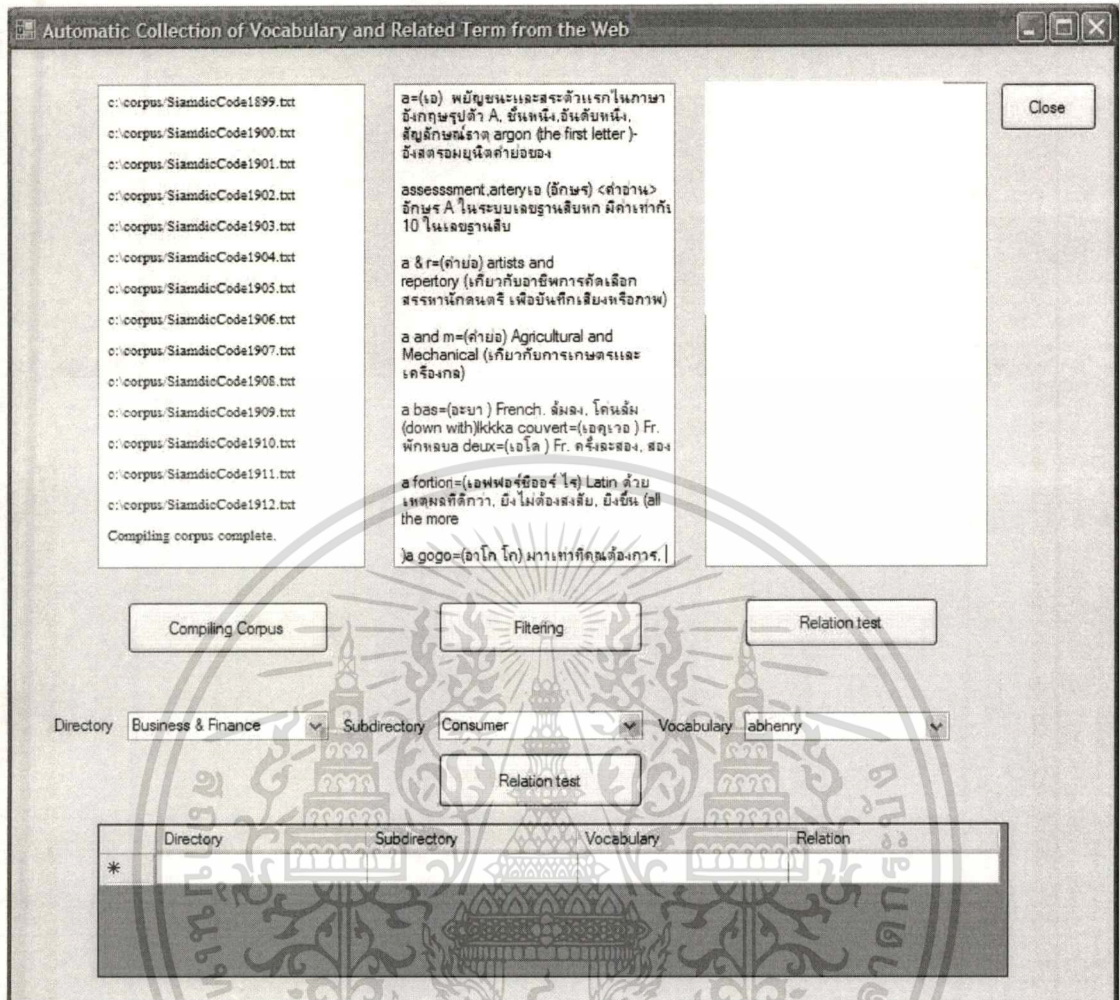
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ ข.9 แสดงผลลัพธ์จากการกดปุ่ม Compiling Corpus

3. กดปุ่ม Filtering ซึ่งเมื่อทำการกดแล้ว การทำงาน Filtering เพื่อดึงคำศัพท์และคำแปลออกมาจาก Corpus ที่เรารวบรวมได้จากการทำ compiling corpus เมื่อการทำงานเสร็จสมบูรณ์จะปรากฏหน้าจอเมนูหลักดังรูป ข.10

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ ข.10 แสดงผลลัพธ์จากการกดปุ่ม Filtering

4. การทำการกดปุ่ม Relation test เพื่อทำการหาความสัมพันธ์ระหว่างคำศัพท์และหมวดหมู่ โดยจะเริ่มทำงานเมื่อมีการกดปุ่ม Relation test เมื่อการทำงานเสร็จสมบูรณ์จะปรากฏหน้าจอเมนูหลักดังรูป ข.11

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Automatic Collection of Vocabulary and Related Term from the Web

c:\corpus/SiamdicCode1899.txt
 c:\corpus/SiamdicCode1900.txt
 c:\corpus/SiamdicCode1901.txt
 c:\corpus/SiamdicCode1902.txt
 c:\corpus/SiamdicCode1903.txt
 c:\corpus/SiamdicCode1904.txt
 c:\corpus/SiamdicCode1905.txt
 c:\corpus/SiamdicCode1906.txt
 c:\corpus/SiamdicCode1907.txt
 c:\corpus/SiamdicCode1908.txt
 c:\corpus/SiamdicCode1909.txt
 c:\corpus/SiamdicCode1910.txt
 c:\corpus/SiamdicCode1911.txt
 c:\corpus/SiamdicCode1912.txt
 Compiling corpus complete.

อ=(เจ) พยัญชนะและสระตัวแรกในภาษาอังกฤษรูปตัว A, ซี, เอ็น, อินดิคทีฟ, สัญลักษณ์แรก aagon (the first letter)- จึงสตรจมยูนิตคำมอชง

assessment, artery เจ (อังกฤษ) <คำอ่าน> อักษร A ในระบยเลขฐานสิบหก มีค่าเท่ากับ 10 ในเลขฐานสิบ

อ & r=(คำย่อ) artists and repertory (เกี่ยวกับอาชีพการคัดเลือกสรรหานักแสดง เพื่อบันเทิงเสียงหรือภาพ)

a and m=(คำย่อ) Agricultural and Mechanical (เกี่ยวกับกาเกษตรและเครื่องกล)

อ bss=(อะนา) French, ส้มลง, โคนฉิม (down with)kkka couvert=(เอตุเออ) Fr. พักหลบ deus=(เอโต) Fr. ครึ่งละสอง, สอง

a fortion=(เอฟฟอจีมอจ ไซ) Latin ตามเหตุผลที่ดีกา, มีงไม่ตองสงสัย, มีงขึ้น (all the more

ja gogo=(อาโก โก) มาเท่ากับกชตองการ.

Business -> yard rope|
 Business -> yard measure
 Business -> yellow journalism
 Business -> yeomans service
 Business -> young man
 Business -> zebec
 Business -> zephyr cloth
 Business -> zeus
 Business -> zip code
 Relation test Complete!!

Close

Compiling Corpus

Filtering

Relation test

Directory Business & Finance Subdirectory Consumer Vocabulary abhenry

Relation test

Directory	Subdirectory	Vocabulary	Relation
*			

รูปที่ ข.11 แสดงผลลัพธ์ที่จากการคปุม Relation test

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ประวัติผู้เขียน

ชื่อผู้เขียน นางสาวอมอร สิริสุภางค์
 สถานที่เกิด จังหวัดกรุงเทพ
 การศึกษา ระดับปริญญาตรี
 วส.บ. (วิทยาการสารสนเทศบัณฑิต)
 สาขาวิชาวิทยาการคอมพิวเตอร์
 มหาวิทยาลัยธรรมศาสตร์

ประสบการณ์การทำงาน Programmer

Samart Info media
 Programmer
 Ictus
 Software testing engineer
 Microsoft(Thailand)



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้