

เครื่องมือท่องเว็บไซต์บนอุปกรณ์เคลื่อนที่สำหรับผู้พิการทางสายตาโดยใช้การตรวจหาเนื้อหาหลัก



บทคัดย่อและแฟ้มข้อมูลฉบับเต็มของวิทยานิพนธ์ตั้งแต่ปีการศึกษา 2554 ที่ให้บริการในคลังปัญญาจุฬาฯ (CUIR)
เป็นแฟ้มข้อมูลของนิสิตเจ้าของวิทยานิพนธ์ ที่ส่งผ่านทางบัณฑิตวิทยาลัย

The abstract and full text of theses from the academic year 2011 in Chulalongkorn University Intellectual Repository (CUIR)
are the thesis authors' files submitted through the University Graduate School.

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรมหาบัณฑิต
สาขาวิชาวิทยาศาสตร์คอมพิวเตอร์ ภาควิชาวิศวกรรมคอมพิวเตอร์
คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย
ปีการศึกษา 2558
ลิขสิทธิ์ของจุฬาลงกรณ์มหาวิทยาลัย

A WEB NAVIGATION TOOL ON MOBILE DEVICE FOR VISUALLY IMPAIRED PERSONS
USING WEB MAIN CONTENT DETECTION

Mr. Krit Bannachaisirisuk



A Thesis Submitted in Partial Fulfillment of the Requirements
for the Degree of Master of Science Program in Computer Science

Department of Computer Engineering

Faculty of Engineering

Chulalongkorn University

Academic Year 2015

Copyright of Chulalongkorn University

หัวข้อวิทยานิพนธ์ เครื่องมือทองเว็บไซด์บนอุปกรณ์เคลื่อนที่สำหรับผู้พิการ
ทางสายตาโดยใช้การตรวจหาเนื้อหาหลัก
โดย นายกฤษฏี บรรณะชัยศิริสุข
สาขาวิชา วิทยาศาสตร์คอมพิวเตอร์
อาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก ผู้ช่วยศาสตราจารย์ ดร. สุกรี สิ้นธุภิญโญ

คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย อนุมัติให้บัณฑิตวิทยานิพนธ์ฉบับนี้เป็นส่วน
หนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรบัณฑิต

.....คณบดีคณะวิศวกรรมศาสตร์
(รองศาสตราจารย์ ดร. สุพจน์ เตชวรสินสกุล)

คณะกรรมการสอบวิทยานิพนธ์

.....ประธานกรรมการ
(ผู้ช่วยศาสตราจารย์ ดร. นัทธี นิภานันท์)

.....อาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก
(ผู้ช่วยศาสตราจารย์ ดร. สุกรี สิ้นธุภิญโญ)

.....กรรมการ
(ผู้ช่วยศาสตราจารย์ ดร. ธีรพงศ์ ชินธเนศ)

.....กรรมการภายนอกมหาวิทยาลัย
(ดร. เด่นดวง ประดับสุวรรณ)

กฤษฎี บรรณัชชัยศิริสุข : เครื่องมือท่องเว็บไซต์บนอุปกรณ์เคลื่อนที่สำหรับผู้พิการทางสายตาโดยใช้การตรวจหาเนื้อหาหลัก (A WEB NAVIGATION TOOL ON MOBILE DEVICE FOR VISUALLY IMPAIRED PERSONS USING WEB MAIN CONTENT DETECTION) อ.ที่ปรึกษาวิทยานิพนธ์หลัก: ผศ. ดร. สุกรี สินธุภิญโญ, 59 หน้า.

เนื่องจากในปัจจุบันเราอยู่ในยุคดิจิทัล ที่ซึ่งเทคโนโลยีกลายเป็นส่วนหนึ่งในชีวิตประจำวัน ข่าวสาร สารระ และความบันเทิงมากมายนั้นสามารถหาได้จากบนอินเทอร์เน็ต อยู่ในรูปแบบของเว็บไซต์ที่มีการเกิดใหม่ขึ้นในทุกๆวัน และมีจำนวนเกือบหนึ่งพันล้านเว็บไซต์ที่ใช้งานอยู่ในขณะนี้ แต่ผู้พิการทางสายตายังคงมีความลำบากในการเข้าถึงเนื้อหาหลักบนเว็บไซต์ งานวิจัยนี้มุ่งที่จะสร้างเครื่องมือเพื่ออำนวยความสะดวกให้กับผู้พิการทางสายตา ให้สามารถเข้าถึงเว็บไซต์ผ่านทางอุปกรณ์พกพาที่มีราคาถูก โดยระบบจะทำหน้าที่ในการรับคำค้นหาจากผู้ใช้ และนำไปค้นหาหัวข้อที่เกี่ยวข้อง พร้อมทั้งสร้างรายการของเว็บไซต์ที่เกี่ยวข้อง จากนั้นเมื่อผู้ใช้สนใจในเว็บไซต์ใด เครื่องมือจะทำการค้นหาเนื้อหาหลักด้วยการนำวิธีคุณลักษณะข้อความแบบตื้น (Shallow text feature) ร่วมกับการใช้คำค้น เพื่อนำไปสร้างตัวคัดแยกด้วยหลักการป่าแบบสุ่ม (Random forest) เมื่อได้ตัวคัดแยก และทำการเรียนรู้จึงจะสามารถนำไปใช้เพื่อทำการคัดแยกเนื้อหาหลักภายในหน้าเว็บ จากนั้นจึงนำเสนอในรูปแบบของตัวอักษรที่ถูกเรียบเรียงใหม่ และทำให้ผู้พิการทางสายตาสามารถเข้าถึงเนื้อหาได้ผ่านทางโปรแกรมอ่านหน้าจอ (Screen reader) ซึ่งทำหน้าที่อ่านข้อความที่ปรากฏให้ผู้พิการทางสายตามีความสามารถในการเข้าถึงเนื้อหาผ่านทางเสียง ด้วยวิธีการนี้จะทำให้ผู้พิการทางสายตามีความสามารถในการเข้าถึงเนื้อหาได้อย่างรวดเร็วมากยิ่งขึ้น อีกทั้งลดขั้นตอนและความซับซ้อนในการเข้าถึงเนื้อหาบนหน้าเว็บไซต์

ภาควิชา วิศวกรรมคอมพิวเตอร์

ลายมือชื่อนิสิต

สาขาวิชา วิทยาศาสตร์คอมพิวเตอร์

ลายมือชื่อ อ.ที่ปรึกษาหลัก

ปีการศึกษา 2558

5671037921 : MAJOR COMPUTER SCIENCE

KEYWORDS: WEB NAVIGATION TOOL / THE VISUALLY IMPAIRED / CONTENT DETECTION

KRIT BANNACHAISIRISUK: A WEB NAVIGATION TOOL ON MOBILE DEVICE FOR VISUALLY IMPAIRED PERSONS USING WEB MAIN CONTENT DETECTION.

ADVISOR: ASST. PROF. DR. SUKREE SINTHUPINYO, 59 pp.

Even though, an endless resource of information is currently available over the internet, but the visually impaired persons are still not able to easily access to content of the website because of disability. In this paper, we introduce a web navigation tool on mobile devices for helping those to be able to get into the main content of website faster. With the help of shallow text features integrated with keywords from the user, a classifier will be constructed by using the random forest method. The classifier will then be applied to remove the boilerplate and extract actual content from the webpage. Screen reader is also required here for reading the extracted main content aloud to the user who cannot navigate using their sight.

The expected outcome from this research is to create a mobile application that improves the accessibility of visually impaired person to get into the actual content faster and easier than existing methodologies.

จุฬาลงกรณ์มหาวิทยาลัย
CHULALONGKORN UNIVERSITY

Department: Computer Engineering Student's Signature

Field of Study: Computer Science Advisor's Signature

Academic Year: 2015

กิตติกรรมประกาศ

วิทยานิพนธ์ฉบับนี้สำเร็จลุล่วงไปได้ด้วยความอนุเคราะห์จากผู้ช่วยศาสตราจารย์ ดร. สุกรี สินธุภิญโญ อาจารย์ที่ปรึกษาวิทยานิพนธ์ ได้สละเวลาให้ความรู้ คำปรึกษา ตรวจสอบและแก้ไขข้อผิดพลาดต่างๆ ตลอดจนการกำกับดูแลและคอยติดตามความก้าวหน้า ทำให้การวิจัยนี้สำเร็จไปได้ด้วยดี ผู้วิจัยขอกราบขอบพระคุณเป็นอย่างสูงไว้ ณ โอกาสนี้

ขอขอบพระคุณผู้ช่วยศาสตราจารย์ ดร. นัทธี นิภานันท์ และผู้ช่วยศาสตราจารย์ ดร. ณัฐพงศ์ ชินธเนศ กรรมการสอบวิทยานิพนธ์ รวมถึง ดร. เต๋นดวง ประดับสุวรรณ กรรมการภายนอกมหาวิทยาลัยที่กรุณาเสียสละเวลา ให้คำแนะนำ ตรวจสอบ และแก้ไขวิทยานิพนธ์ฉบับนี้

ขอขอบพระคุณบิดา มารดา และญาติพี่น้องที่ให้การสนับสนุนและเป็นกำลังใจที่ดีให้เสมอมาและสนับสนุนด้านทุนทรัพย์ในการศึกษารวมไปถึงทุกท่านที่มีส่วนช่วยเหลือในการทำวิทยานิพนธ์ครั้งนี้ ซึ่งมีได้กล่าวนามในที่นี้

ท้ายที่สุด ผู้วิจัยขอขอบพระคุณเพื่อนๆ ทุกคน ที่คอยติดตามและให้กำลังใจ รวมถึงท่านอื่นๆ ที่มีได้กล่าวลงนามไว้ ณ ที่นี้ที่มีส่วนทำให้วิทยานิพนธ์สำเร็จลุล่วงไปได้ด้วยดีผู้วิจัยหวังเป็นอย่างยิ่งว่าวิทยานิพนธ์ฉบับนี้จะเป็นประโยชน์บ้างไม่มากก็น้อยสำหรับผู้สนใจจะศึกษารายละเอียดต่อไป

สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	ง
บทคัดย่อภาษาอังกฤษ.....	จ
กิตติกรรมประกาศ.....	ฉ
สารบัญ.....	ช
บทที่ 1 บทนำ	1
1.1 ความเป็นมาและความสำคัญของปัญหา	1
1.2 วัตถุประสงค์ของการวิจัย.....	2
1.3 ขอบเขตงานวิจัย.....	2
1.4 ขั้นตอนและวิธีการดำเนินการวิจัย	3
1.5 ประโยชน์ที่คาดว่าจะได้รับ.....	3
1.6 ผลงานตีพิมพ์จากวิทยานิพนธ์.....	4
1.7 โครงสร้างของเนื้อหาในวิทยานิพนธ์.....	4
บทที่ 2 ทฤษฎีและงานวิจัยที่เกี่ยวข้อง	5
2.1. ทฤษฎีที่เกี่ยวข้อง	5
2.1.1. การแปลงเอกสารเว็บไซต์ (HTML Parsing).....	5
2.1.2. คุณสมบัติแบบตื้น (Shallow text feature)	6
2.1.3. วิธีการป่าแบบสุ่ม (Random forest)	6
2.1.4. เว็บบริการ (Web Service)	8
2.2. งานวิจัยที่เกี่ยวข้อง.....	8
2.2.1. Application of Content Adaptation in Web Accessibility for the Blind [4].....	8
2.2.2. Boilerplate Detection using Shallow Text Features [2].....	10

2.2.3. Extracting news text from web pages: an application for the visually impaired [5].....	13
บทที่ 3 แนวคิดและวิธีดำเนินการวิจัย	15
3.1. การส่งคำร้องขอรายการเว็บไซต์จากคำค้นหา	16
3.2. การแปลงข้อมูลเอกสารเว็บไซต์และแสดงผล	17
3.3. การสังเคราะห์ข้อมูลเบื้องต้นและการสร้างข้อมูลสำหรับการเรียนรู้จากหน้าเว็บไซต์	19
3.4. การติดต่อเว็บบริการและค้นหาเนื้อหาหลัก.....	22
3.5. การแสดงผลสำหรับเครื่องมืออ่านหน้าจอ	23
บทที่ 4 การพัฒนาเครื่องมือ.....	25
4.1. ความต้องการเชิงฟังก์ชัน.....	25
4.2. การวิเคราะห์ความต้องการและแผนภาพฟังก์ชันงานของระบบ.....	25
4.3. สภาพแวดล้อมที่ใช้ในการพัฒนาเครื่องมือสนับสนุน.....	27
4.4. ขั้นตอนการทำงานของเครื่องมือ	28
4.5 ขั้นตอนการจัดการระบบเว็บบริการ	32
บทที่ 5 การทดสอบและการวิเคราะห์ผล.....	37
5.1. วัตถุประสงค์ของการทดสอบ.....	37
5.2. การทดสอบระบบ.....	37
5.3. สรุปผลการทดสอบ	51
บทที่ 6 สรุปผลการวิจัย.....	52
6.1. สรุปผลการวิจัย	52
6.2. ข้อจำกัดของงานวิจัย.....	53
6.3. งานวิจัยในอนาคต.....	53
รายการอ้างอิง.....	54

ภาคผนวก ก. การติดตั้งส่วนเสริมโคโคพอด (Cocoa Pod), อลาโมไฟร์ (Alamofire) และตัว อ่านเอกสารเว็บไซต์ (HTML Reader).....	56
ภาคผนวก ข. การติดตั้งส่วนต่อประสานเว็บบริการ (bigml-swift).....	57
ประวัติผู้เขียนวิทยานิพนธ์	59



บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

ในยุคที่มีอินเทอร์เน็ตเป็นส่วนหนึ่งในชีวิตประจำวัน มีผู้ให้บริการเว็บไซต์จำนวนมากนับพันล้านเว็บไซต์ในขณะนี้ (ที่มา: <http://www.internetlivestats.com/total-number-of-websites>) ซึ่งประกอบไปด้วยเนื้อหาและข้อมูลมากมายที่ถูกเขียนเพิ่มขึ้นทุกๆวินาทีทั้ง ข่าวสาร สารระ และความบันเทิง แต่กลับมีผู้ใช้ที่ไม่สามารถเข้าถึงเนื้อหาเหล่านั้นได้อย่างเต็มที่ อันเนื่องมาจากความไม่สมบูรณ์ทางกายภาพที่ไม่สามารถมองเห็นได้เหมือนผู้ใช้ปกติทั่วไป มีเทคโนโลยีช่วยเหลือ (Assistive Technology) ที่เรียกว่า เครื่องมืออ่านหน้าจอ (Screen Reader) ซึ่งถูกพัฒนาขึ้นมาเพื่อใช้ในการแปลงตัวอักษรที่อยู่บนหน้าจอออกมาเป็นรูปแบบของเสียงสังเคราะห์ทำให้ผู้ใช้ที่มีความพิการทางสายตาสามารถเข้าถึงข้อมูลเหล่านั้นได้ โดยการพัฒนาหน้าเว็บไซต์ที่จะทำให้เครื่องมืออ่านหน้าจอมีประสิทธิภาพสูงสุด จะต้องมีการพัฒนาตามหลักแนวทางที่ชื่อว่า Web Content Accessibility Guidelines 2.0 (WCAG 2.0) ถูกกำหนดขึ้นโดยองค์กร World Wide Web Consortium (W3C) (www.w3.org/) แนวทางการพัฒนานี้จะทำให้หน้าเว็บไซต์สามารถรองรับการเข้าถึงจากเครื่องมืออ่านหน้าจอได้ดีขึ้น เนื่องจากการวางโครงสร้างของเอกสารเว็บไซต์ (HTML Document Object Model) ที่แบ่งออกเป็นระดับแบบโครงสร้างต้นไม้ (DOM Tree) ทำให้ผู้พิการสามารถเข้าไปยังจุดที่ต้องการได้อย่างง่ายดายและรวดเร็วมากยิ่งขึ้น แต่อย่างไรก็ตามมีเว็บไซต์เพียงจำนวนเล็กน้อยเท่านั้นที่ถูกพัฒนาตามแนวทางของ WCAG 2.0 [1] ทำให้เครื่องมืออ่านหน้าจอยังคงขาดความสามารถที่จะทำให้ผู้ใช้สามารถเข้าถึงเนื้อหาหลักบนหน้าเว็บไซต์ได้อย่างรวดเร็ว การที่ผู้พิการจะต้องใช้เครื่องมืออ่านหน้าจอและฟังในหลายๆหัวข้อ นั้น มีผลทำให้ต้องเสียเวลาในการค้นหาเนื้อหาด้วยตัวเอง บนหน้าเว็บไซต์นั้นๆ ซึ่งอาจจะประกอบไปด้วยส่วนประกอบมากมายที่ไม่ได้เกี่ยวข้องกับเนื้อหาที่ผู้พิการต้องการ เช่น เมนู, โฆษณา หรือ บทความที่เกี่ยวข้องอื่นๆ ดังนั้นการนำวิธีการค้นหาและจำแนกเนื้อหาหลักของหน้าเว็บไซต์จึงเป็นสิ่งสำคัญที่จะนำมาใช้เพื่อค้นหาและเรียงเรียงข้อความมากมายในหน้าเว็บไซต์นั้นๆออกมาเป็นรูปแบบใหม่ที่จะทำให้ผู้พิการทางสายตาสามารถเข้าถึงตัวเนื้อหาหลักได้อย่างรวดเร็ว

จึงเกิดเป็นแรงบันดาลใจในงานวิจัยนี้เพื่อที่จะนำหลักการต่างๆที่ได้ศึกษามา นำไปใช้ในการสร้างเครื่องมือเพื่ออำนวยความสะดวกให้กับผู้พิการทางสายตา เพื่อให้เข้าถึงเนื้อหาเว็บไซต์ที่ต้องการผ่านทางเครื่องมือสื่อสารแบบพกพา ผ่านการป้อนคำค้นหาที่ได้มาจากความสนใจของผู้ใช้ มาแสดงเป็นรายการเว็บไซต์ที่มีความเชื่อมโยงกับคำค้นหาแล้ว จากนั้นเมื่อผู้ใช้ได้ทำการเลือกที่อยู่ของเว็บไซต์

(Web Address) จากนั้นทำการคัดแยกเนื้อหาที่แท้จริงจากเว็บไซต์นั้นๆ ออกมาโดยวิธีการที่จะทำการคัดแยกเนื้อหาหน้าเว็บนั้น สิ่งที่จะต้องเข้าใจอย่างหนึ่งคือโครงสร้างของเว็บไซต์ที่ถูกสร้างขึ้นจากภาษา HTML (HyperText Markup Language) ไม่ได้ประกอบไปด้วยเพียงส่วนที่เป็นเนื้อหาเท่านั้น ยังประกอบไปด้วยส่วนที่มองไม่เห็นซึ่งมีลักษณะที่เป็นภาษาโปรแกรมต่างๆ เช่น JavaScript, Cascading Style Sheet (CSS) เป็นต้น โดยส่วนต่างๆ เหล่านี้ไม่ได้มีผลสำคัญต่อผู้พิการทางสายตาอันเนื่องมาจาก ผู้ใช้นั้นจะไม่สามารถเข้าถึงส่วนที่เป็นการแสดงผลเหล่านั้นได้ ดังนั้นส่วนเหล่านี้จึงควรที่จะถูกคัดแยกออกจากส่วนที่เป็นเนื้อหาหลักของหน้าเว็บ จากการวิจัยค้นพบว่า กระบวนการที่เรียกว่าคุณลักษณะแบบตื้น (Shallow feature text) เป็นวิธีการที่เหมาะสม รวดเร็ว และแม่นยำที่สุดในขณะนี้ ที่จะนำไปใช้เพื่อจำแนกเนื้อหาหลักในหน้าเว็บ งานวิจัยนี้ได้นำเสนอเครื่องมือในรูปแบบของแอปพลิเคชันบนเครื่องมือสื่อสารที่จะทำให้ผู้พิการทางสายตามีความสามารถในการเข้าถึงเนื้อหาหรือข้อมูลต่างๆ บนหน้าเว็บไซต์ได้อย่างรวดเร็ว และมีประสิทธิภาพมากขึ้นกว่าวิธีการในปัจจุบัน

1.2 วัตถุประสงค์ของการวิจัย

นำเสนอวิธีการและพัฒนาเครื่องมือสำหรับช่วยเหลือผู้พิการทางสายตาในการเข้าถึงเนื้อหาเว็บไซต์ได้ผ่านเครื่องมือสื่อสารเคลื่อนที่ระบบไอโอเอส โดยใช้การค้นหาเนื้อหาหลักของหน้าเว็บไซต์ด้วยวิธีการป่าแบบสุ่ม และนำมาเปลี่ยนแปลงการนำเสนอในรูปแบบใหม่ ที่จะทำให้การเข้าถึงจากผู้พิการทางสายตาเป็นไปอย่างรวดเร็วมากขึ้น

1.3 ขอบเขตงานวิจัย

- 1) ผู้ใช้จะต้องทำการป้อนคำค้นหาให้กับเครื่องมือด้วยตัวเอง เพื่อทำการค้นหารายการเว็บไซต์ที่มีความเกี่ยวข้องจากคำเหล่านั้น
- 2) ส่วนของการค้นหาเว็บไซต์ที่เกี่ยวข้องกับคำค้นหา จะได้ผลลัพธ์มาจาก Google.com ด้วยการแปลงเอกสารเว็บไซต์ (HTML Parsing)
- 3) ส่วนของการวิเคราะห์เพื่อค้นหาเนื้อหาหลักจะได้มาจากเว็บบริการ จากผู้ให้บริการ BigML.com ด้วยวิธีการป่าแบบสุ่ม (Random Forest)
- 4) เครื่องมือในงานวิจัยนี้จะอาศัยเครื่องมืออ่านหน้าจอบนระบบไอโอเอสที่ชื่อว่า VoiceOver ในการอ่านส่วนประกอบต่างๆบนหน้าจอ รวมถึงเนื้อหาหลักที่ได้หลังถูกจำแนกจากเว็บบริการ

1.4 ขั้นตอนและวิธีการดำเนินการวิจัย

- 1) ศึกษาและทำความเข้าใจวิธีการเข้าถึงเว็บไซต์ของผู้พิการทางสายตา
- 2) ศึกษาและทำความเข้าใจถึงปัญหาในการเข้าถึงเนื้อหาของผู้พิการทางสายตา
- 3) ศึกษาและทำความเข้าใจถึงโปรแกรมอ่านหน้าจอ
- 4) ศึกษาและทำความเข้าใจถึงวิธีการที่มีอยู่ในปัจจุบัน
- 5) ศึกษาและทำความเข้าใจถึงการประยุกต์ใช้ทฤษฎีการเรียนรู้ของเครื่อง (Machine Learning)
- 6) ออกแบบโครงสร้างและโปรแกรม
- 7) พัฒนาระบบ
- 8) ทดสอบและประเมินผลงานวิจัย
- 9) สรุปผลงานวิจัย และนำผลที่ได้ไปหาจุดบกพร่องพร้อมปรับปรุงแก้ไขให้บรรลุวัตถุประสงค์ของงานวิจัย
- 10) ตีพิมพ์ผลงานทางวิชาการ
- 11) จัดทำวิทยานิพนธ์

1.5 ประโยชน์ที่คาดว่าจะได้รับ

ได้ศึกษาถึงวิธีการและสร้างเครื่องมือที่อำนวยความสะดวกในการท่องเว็บสำหรับผู้พิการทางสายตา ด้วยวิธีการนำวิชาการเรียนรู้ของเครื่อง (Machine Learning) มาใช้สร้างเครื่องมือในการค้นหาเนื้อหาหลักของหน้าเว็บไซต์ และนำเสนอในรูปแบบใหม่เพื่อทำให้การเข้าถึงเนื้อหาเป็นไปได้อย่างรวดเร็วมากยิ่งขึ้น นอกจากนี้ยังได้คำนึงถึงความสำคัญของผู้พิการที่ต้องการเข้าถึงข้อมูลเช่นเดียวกับผู้คนทั่วไป เพียงแต่ขาดโอกาสที่จะได้รับและผู้สนับสนุนที่เพียงพอจากสังคมในปัจจุบัน การนำเทคโนโลยีมาพัฒนาเครื่องมือจึงเป็นคำตอบที่จะสามารถเข้าไปช่วยอำนวยความสะดวกแก่ผู้พิการทางสายตา

1.6 ผลงานตีพิมพ์จากวิทยานิพนธ์

ส่วนหนึ่งของวิทยานิพนธ์นี้ได้รับการตีพิมพ์ในรายงานสืบเนื่องจากการประชุมวิชาการระดับชาติเรื่อง “A Web Navigation Tool on Mobile Device for Visually Impaired Persons Using Web Main Content Detection, Proceedings of 12th National Conference on Computing and Information Technology (NCCIT 2016), July 7-8, 2016, Khon Kaen, Thailand” และได้รับรางวัล Best Paper Award ในหัวข้อ Information Technology and System Engineering.

1.7 โครงสร้างของเนื้อหาในวิทยานิพนธ์

เนื้อหาของวิทยานิพนธ์ฉบับนี้แบ่งออกเป็น 6 บทด้วยกันคือ บทที่ 1 อธิบายถึงที่มาและความสำคัญของปัญหา รวมถึงขอบเขตและประโยชน์ของงานวิจัย บทที่ 2 อธิบายถึงทฤษฎีที่เกี่ยวข้องและงานวิจัยที่เกี่ยวข้อง บทที่ 3 อธิบายถึงแนวคิดและวิธีการดำเนินการวิจัยในการสร้างเครื่องมือเพื่อสร้างรายการเว็บไซต์จากคำค้นหา และการใช้เว็บบริการในการค้นหาเนื้อหาหลักบนหน้าเว็บไซต์ บทที่ 4 อธิบายถึงวิธีการพัฒนาเครื่องมือสนับสนุนแนวคิดของงานวิจัย บทที่ 5 อธิบายวิธีการทดสอบและการวิเคราะห์ผลและในบทสุดท้ายจะสรุปงานวิจัยทั้งหมด รวมถึงงานวิจัยในอนาคต



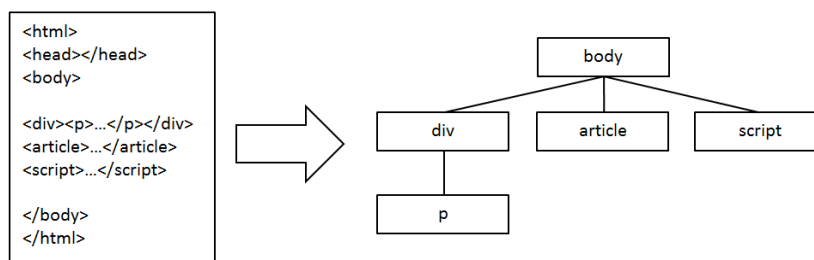
บทที่ 2

ทฤษฎีและงานวิจัยที่เกี่ยวข้อง

2.1. ทฤษฎีที่เกี่ยวข้อง

2.1.1. การแปลงเอกสารเว็บไซต์ (HTML Parsing)

การแปลงเอกสารเว็บไซต์คือ การนำเอกสารเว็บไซต์ที่อยู่ในภาษา HTML (Hypertext Markup Language) มาทำการแปลงให้อยู่ในรูปแบบใหม่ โดยในขั้นแรกคือการส่งคำร้องขอไปยังที่อยู่ของเว็บไซต์จากนั้นผู้ให้บริการ หรือ เซิร์ฟเวอร์ (Server) จะตอบกลับในรูปแบบเอกสารเว็บไซต์ (HTML Document) ที่ประกอบไปด้วยส่วนประกอบต่างๆในหน้าเว็บไซต์นั้น ซึ่งงานวิจัยนี้จะทำการแปลงเอกสารเว็บไซต์ที่ได้จากรูปแบบตัวอักษรให้อยู่ในรูปแบบใหม่เป็นลักษณะของโครงสร้างต้นไม้ ที่เรียกว่า เอกสารเว็บไซต์แบบอ็อบเจกต์โมเดลต้นไม้ (HTML Document Object Model Tree) เพื่อให้เครื่องมือสามารถเข้าถึงแต่ละส่วนประกอบได้สะดวกมากขึ้น



รูปที่ 1 ภาพรวมการแปลงเอกสารเว็บไซต์เป็นแบบอ็อบเจกต์โมเดลต้นไม้

ในรูปที่ 1 การแปลงเอกสารเว็บไซต์จะนำส่วนประกอบแต่ละส่วนมาจัดอยู่ในรูปแบบโครงสร้างต้นไม้ โดยแต่ละโนด (Node) จะมีความสัมพันธ์กับเอกสารเว็บไซต์จากความสัมพันธ์ของชื่อแท็กข้อมูล (Tag) และนำแต่ละส่วนประกอบมาสร้างเป็นแบบจำลองข้อมูลเชิงลำดับชั้น (Hierarchical Data Model)

การแปลงเอกสารเว็บไซต์ไปเป็นโมเดลต้นไม้จะช่วยทำให้ผู้พัฒนาสามารถนำโมเดลที่ได้ไปสู่กระบวนการอื่นๆ ไม่ว่าจะเป็นการค้นหา จัดเรียงส่วนประกอบ หรือการคัดกรองส่วนที่ไม่สามารถเป็นเนื้อหาได้จากการพิจารณาจากชื่อแท็กเช่น รูปภาพ (img), ชุดคำสั่ง (script) หรือ ข้อมูลเมตา (meta) ซึ่งเป็นเพียงส่วนประกอบของหน้าเว็บไซต์ เป็นต้น

2.1.2. คุณสมบัติแบบตื้น (Shallow text feature)

การค้นหาเนื้อหาหลักบนเว็บไซต์จากการนำข้อมูลที่ได้จากการนำคุณลักษณะเบื้องต้นของแต่ละองค์ประกอบมาประมวลผล และนำไปใช้ในการค้นหาส่วนที่เป็นเนื้อหาหลักได้จากคุณสมบัติแบบตื้น (Shallow text feature) [2] ร่วมกับวิธีการเรียนรู้ของเครื่อง (Machine Learning) ซึ่งประกอบไปด้วยคุณสมบัติดังนี้

a. ความหนาแน่นของตัวอักษร (Text Density)

คำนวณจากการนำจำนวนตัวอักษรทั้งหมดในบล็อก b_x มาทำการหารด้วยตัวอักษรในหนึ่งบรรทัดซึ่งจะถูกกำหนดด้วยผู้พัฒนา โดยมาตรฐานสำหรับภาษาอังกฤษจะอยู่ที่ 80 - 90 ตัวอักษรต่อบรรทัด

$$q(b_x) = \frac{\text{Number of token in } b_x}{\text{Number of token per line}}$$

b. ความหนาแน่นของลิงก์ (Link Density หรือ Anchor Percentage)

คำนวณจากการนำจำนวนตัวอักษรที่อยู่ภายใต้แท็ก A (ตัวหนังสือที่จะถูกใช้เพื่อเป็นลิงก์) ในบล็อก b_x มาทำการหารด้วยจำนวนตัวอักษรทั้งหมดในบล็อกเดียวกันนั้น

$$\alpha(b_x) = \frac{\text{Number of token under A tag in } b_x}{\text{Number of token under any tags in } b_x}$$

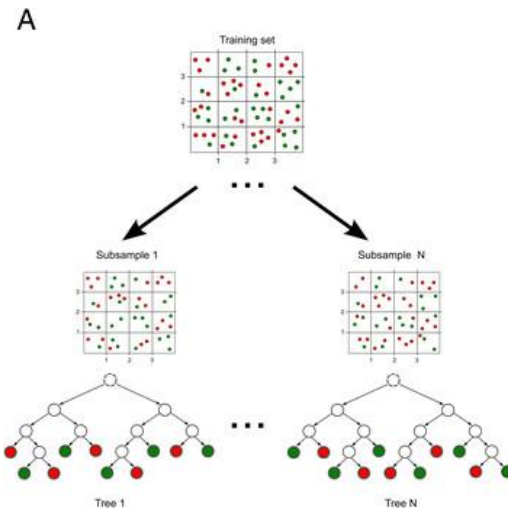
c. น้ำหนักของคำค้นหา (Keywords weight)

ตามแนวคิดของผู้วิจัย ต้องการนำคำค้นหามาเป็นทดสอบว่าเป็นส่วนหนึ่งในการค้นหาเนื้อหาหลักหรือไม่ โดยจะถูกคำนวณจากการปรากฏคำค้นหาภายในบล็อก b_x และในกรณีที่มีการใช้คำค้นหามากกว่าหนึ่งคำ จะให้น้ำหนักในคำค้นหาคำแรกมากที่สุด และไล่เรียงลงมาตามลำดับ

$$\beta(b_x) = \text{Keyword Weight} \times \text{Occurances}$$

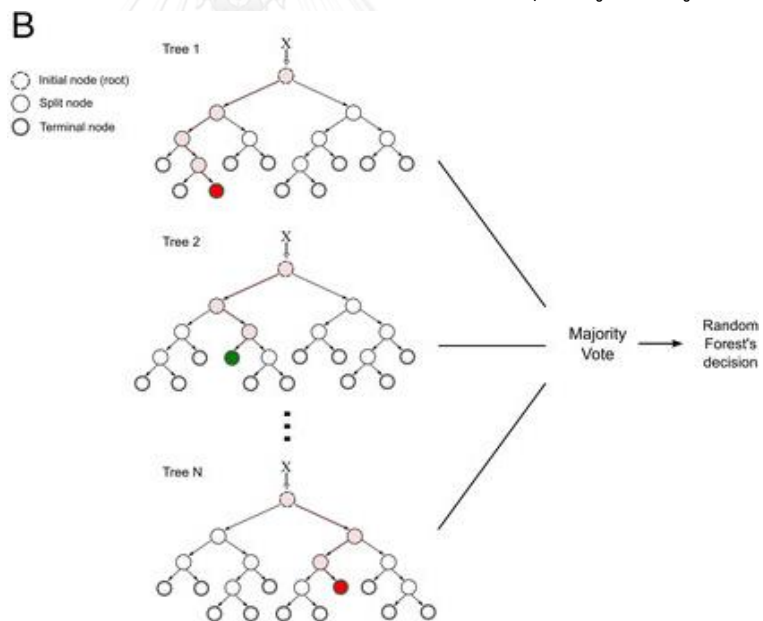
2.1.3. วิธีการป่าแบบสุ่ม (Random forest)

วิธีการป่าแบบสุ่มคือหนึ่งในเทคนิคการสร้างตัวจำแนกแบบชุดที่ถูกคิดค้นขึ้นโดย Breiman [3] อาศัยการสร้างโมเดลสำหรับจำแนก จำนวนหลายๆโมเดลมาช่วยในการค้นหาคำตอบ ซึ่งวิธีการป่าแบบสุ่มจะมีเทคนิคในการสร้างต้นไม้ตัดสินใจหลายๆต้น และทำการสุ่มเลือกแอตทริบิวต์ต่างๆ เป็นหลายๆชุด และนำไปสร้างโมเดลด้วยการต้นไม้ตัดสินใจจำนวนหลายๆต้น โดยการทำนายคำตอบจะกระทำโดยการนำผลโหวตจากต้นไม้ตัดสินใจมาวิเคราะห์เพื่อหาว่าคำตอบเป็นสิ่งที่ใดมากที่สุด (Majority Vote) สุดท้ายจึงสรุปออกมาได้ว่าผลลัพธ์ในการจำแนกคืออะไร



รูปภาพที่ 2 แสดงแนวคิดของวิธีการป่าแบบสุ่มสำหรับการแยกปัญหา 2 คลาส [3]

จากรูปที่ 2 เป็นการแสดงการสร้างการตัดสินใจจากชุดข้อมูลสำหรับการเรียนรู้ (Training set) โดยประกอบไปด้วยสองปัญหาคือ สีเขียวและสีแดง ซึ่งข้อมูลแบบย่อย (Subsample) แต่ละอันสามารถนำไปสร้างต้นไม้ตัดสินใจได้แบบแตกต่างกันออกไปจากการนำข้อมูลแบบสุ่ม (Random bootstrap sample) ออกมาจากชุดข้อมูลเรียนรู้เดิม



รูปที่ 3 แสดงวิธีการโหวตเพื่อหาคำตอบของปัญหาจากการใช้วิธีป่าแบบสุ่ม [3]

จากรูปที่ 3 หลังจากการสร้างชุดต้นไม้ตัดสินใจ จำนวน N ต้นจากการนำชุดข้อมูลที่ได้จากการนำชุดข้อมูลเรียนรู้มาทำการสุ่ม เมื่อนำชุดข้อมูลทดสอบที่ต้องการหาคำตอบเข้าไปสู่ตัวคัดแยก ต้นไม้ตัดสินใจแต่ละต้นจะทำการโหวตเพื่อให้ได้คลาสดำตอบ

2.1.4. เว็บบริการ (Web Service)

เว็บบริการคือการระบบซอฟต์แวร์ที่ออกแบบมา เพื่อช่วยในการให้บริการ แลกเปลี่ยนข้อมูลระหว่างเครื่องมือกับเว็บที่ให้บริการผ่านทางระบบอินเทอร์เน็ต ซึ่งมีส่วนต่อประสานที่มีภาษาที่ใช้อธิบายถึงคุณลักษณะของเว็บบริการนั้นๆ เรียกว่า Web Service Description Language (WSDL) ที่ถูกควบคุมด้วยหน่วยงาน W3C (World Wide Web Consortium) โดยในระบบที่ต้องการจะทำการติดต่อกับเว็บบริการ จะทำการติดต่อผ่านการ ใช้โปรโตคอลพื้นฐานที่เรียกว่า Simple Object Access Protocol (SOAP) โดยติดต่อกัน ผ่านการใช้ภาษา Extensible Markup Language (XML) ที่เป็นภาษามาตรฐานในการ ติดต่อกันระหว่างระบบต่างๆ แต่ต่อมาพบว่า SOAP ที่ตั้งอยู่บนภาษา XML มีความลำบากใน การใช้ในภาษา Javascript อันเนื่องมาจากวิธีในการคัดแยกส่วนที่เป็นเนื้อหาภายใต้แต่ละ แท็กข้อมูลจะต้องอาศัยวิธีการอันซับซ้อน จึงทำให้มีการพัฒนาให้มีการติดต่อแบบใหม่ เรียกว่า Representational State Transfer (REST) ที่มีความซับซ้อนในการเก็บข้อมูล น้อยลง เช่น Comma Separated Value (CSV), JavaScript Object Notation (JSON) หรือ Really Simple Syndication (RSS) ทำให้มีความง่ายต่อการแปลงข้อมูลไปเป็นตัวแปร ในทุกๆภาษามากขึ้น แต่จะต้องอาศัยช่องทางการติดต่อแบบ HyperText Transfer Protocol (HTTP) ร่วมกับการส่งคำร้องขอผ่านทางที่อยู่เว็บไซต์

2.2. งานวิจัยที่เกี่ยวข้อง

2.2.1. Application of Content Adaptation in Web Accessibility for the Blind [4]

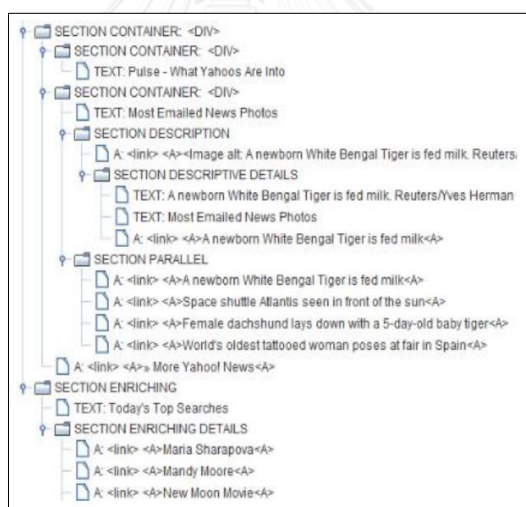
งานวิจัยชิ้นนี้ได้กล่าวถึงปัญหาในการใช้เว็บไซต์ของคนพิการทางสายตา ใน เทคโนโลยีปัจจุบันคนพิการทางสายตาจะสามารถท่องเว็บไซต์ได้ผ่านการใช้โปรแกรมอ่าน หน้าจอ (Screen readers) เช่น Microsoft Narrator หรือ Mac VoiceOver แต่โปรแกรม อ่านหน้าจอก็ยังมีขีดจำกัดในการใช้เพื่อท่องเว็บไซต์ เนื่องจากโปรแกรมเหล่านี้จะอ่านใน ลักษณะจากบนลงล่าง ซ้ายไปขวา หรือขวาไปซ้าย ซึ่งเป็นลักษณะของการอ่านแบบ ตามลำดับ (Sequential) ถ้าหากเนื้อหาที่ต้องการนั้นอยู่ห่างจากส่วนบนมากๆ จะทำให้การ เข้าถึงเนื้อหาเหล่านั้นเป็นไปได้ยาก และใช้เวลานาน

ซึ่งมีการพัฒนาโปรแกรมเพื่อใช้อ่านหน้าจอที่ทันสมัยมากขึ้นชื่อว่า JAWS โปรแกรม นี้มีความสามารถที่จะเข้าถึงข้อมูลได้รวดเร็วมากขึ้น เนื่องจากมีการวิเคราะห์โครงสร้างของ

เว็บไซต์ทำให้ข้ามไปยังจุดต่างๆ ตามที่ต้องการได้ แต่ก็ยังต้องใช้เวลานาน เพราะว่าการอ่านยังคงเป็นลักษณะแบบตามลำดับอยู่ดี เพียงแต่มีความสามารถที่จะข้ามไปยังหัวข้อต่างๆได้

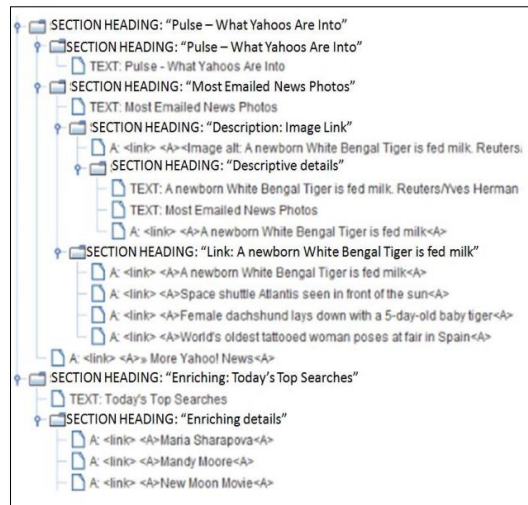
งานวิจัยชิ้นนี้จึงได้เสนอวิธีการที่จะจัดเรียงลำดับโครงสร้างของเว็บไซต์ใหม่ โดยแบ่งกลุ่มเนื้อหาเป็นหัวข้อต่างๆตามความสอดคล้องกัน จากนั้นจึงกำหนดลำดับของหัวข้อและแทรกข้อความเพื่ออธิบายถึงลักษณะของหัวข้อนั้นๆ ซึ่งจะทำให้เกิดเป็นลำดับของข้อมูลแบบมีหัวข้อ จะทำให้ผู้พิการทางสายตาสามารถเข้าถึงข้อมูลที่ต้องการ ผ่านการอ่านหัวข้อแทนที่จะต้องอ่านทั้งหมด ทำให้สามารถเข้าถึงเนื้อหาได้อย่างรวดเร็วมากขึ้น

วิธีการในการจัดเรียงโครงสร้างของหน้าเว็บไซต์ใหม่คือการใช้วิธีที่เรียกว่า วิศวกรรมย้อนรอย (Reverse Engineering) ซึ่งมีขั้นตอนดังนี้ ขั้นแรกจะทำการอ่านเอกสารเว็บไซต์และแยกส่วนประกอบโดยพิจารณาจากโครงสร้างเว็บที่อยู่ในภาษา HTML ในส่วนที่เป็นส่วนประกอบทางตรรกะเช่น ตัวอักษร รูปภาพ ช่องรับอินพุต ปุ่ม ฯลฯ จากนั้นพิจารณาความสัมพันธ์ของแต่ละส่วนว่าจะทำการรวมหรือแยกส่วนไหนออกจากกันแล้วทำการสร้างโมเดลที่เรียกว่า Semantic DOM tree ขึ้นมาจากรูปร่างความสัมพันธ์ของส่วนประกอบดังภาพที่ 4

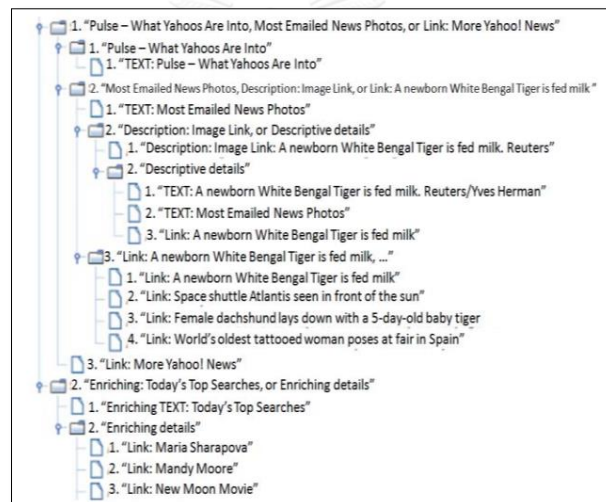


ภาพที่ 4 โมเดลผลลัพธ์จากการแยกส่วนประกอบของหน้าเว็บไซต์ [4]

หลังจากสร้างโมเดลที่ได้จากการแยกรูปประกอบของหน้าเว็บแล้วต่อไปจะทำการกำหนดชื่อหัวข้อให้กับแต่ละกลุ่มโดยมีเงื่อนไขว่าถ้าหากกลุ่มใดมีสมาชิกเพียงหนึ่งตัวก็ใช้ชื่อนั้นเป็นชื่อของหัวข้อกลุ่มนั้นๆ แต่ถ้าหากมีสมาชิกหลายตัว ก็ทำการพิจารณาจากประเภทของสมาชิก และนำไปรวมกันเพื่อสร้างเป็นชื่อกลุ่ม (ดังภาพที่ 5 และภาพที่ 6) ซึ่งเป็นการจัดกลุ่มแบบมีความเกี่ยวข้องกันทางเหตุผล ก็จะทำให้การเข้าถึงข้อมูลเป็นไปอย่างรวดเร็วมากขึ้น



ภาพที่ 5 โมเดลแสดงการใช้ชื่อหัวข้อจากสมาชิกตัวแรก [4]



ภาพที่ 6 โมเดลแสดงการใช้ชื่อหัวข้อจากการพิจารณาสมาชิกทุกตัว [4]

ผู้วิจัยได้สรุปว่า เมื่อมีการจัดกลุ่มของข้อมูลบนหน้าเว็บโดยใช้ความสัมพันธ์กันทางตรรกะจะทำให้ผู้พิการทางสายตาสสามารถเข้าถึงข้อมูลอันซับซ้อนของหน้าเว็บได้อย่างง่ายดาย เนื่องมาจากมีการจัดโครงสร้างใหม่ การเข้าถึงข้อมูลจึงมีลักษณะเป็นลำดับต้นไม้ โดยสามารถหาข้อมูลจากหัวข้อที่เกี่ยวข้องกัน และสามารถข้ามไปยังจุดที่ต้องการผ่านทางหัวข้อได้โดยไม่จำเป็นต้องเข้าถึงส่วนที่เป็นเนื้อหาที่มีปริมาณตัวอักษรมากๆ ได้

2.2.2. Boilerplate Detection using Shallow Text Features [2]

การท่องเว็บของผู้พิการนั้นโดยปกติต้องอาศัยโปรแกรมอ่านหน้าจอ ซึ่งโปรแกรมจะทำหน้าที่อ่านตัวหนังสือที่อยู่บนหน้าเว็บนั้นๆ แต่เนื่องด้วยหน้าเว็บไซต์นั้นอาจจะประกอบไปด้วยข้อความที่ไม่ใช่เนื้อหาหลัก มีลักษณะที่เรียกว่าเป็นแท่งหลอม (Boilerplate) เช่น เมนู

ข้อความเกี่ยวกับผู้จัดทำ ข้อความโฆษณา หรือข้อความอื่นๆที่ไม่ได้เกี่ยวข้องกับเนื้อหา งานวิจัยขึ้นนี้จึงได้คิดค้นวิธีการสังเคราะห์หาข้อความที่มีความเป็นไปได้ว่าจะไม่ใช่เนื้อหาของหน้านั้น ด้วยวิธีการที่เรียกว่า คุณลักษณะแบบตื้น (Shallow Text Features) ซึ่งเป็นวิธีแบบฮิวริสติก (Heuristic) กล่าวคือใช้การศึกษาจากข้อมูลและวิเคราะห์จากการคำนวณ ร่วมกันการใช้อัลกอริทึมต้นไม้ตัดสินใจ โดยงานวิจัยนี้ได้เน้นการค้นหาข้อความเหล่านั้นเพื่อนำไปประมวลผลให้ได้ส่วนที่เป็นเนื้อหาแท้จริงบนเว็บไซต์ และได้ทำการทดลองและฝึกฝน ฟังก์ชันจากกลุ่มของข่าวที่มีมากถึง 254,000 หน้าจาก 408 เว็บไซต์ทั่วโลกด้วยระบบการค้นหาข่าวสารของกูเกิล (Google News Search Engine)

ในงานวิจัยนี้ได้คิดค้นวิธีค้นหาส่วนที่ไม่ต้องการเพื่อที่จะตัดออกจากการวิเคราะห์ เนื้อหา โดยที่คำนวณจากจำนวนของคำ และจำนวนประโยคโดยใช้การตัดคำด้วยช่องว่างเพื่อหาความหนาแน่นของตัวอักษรในส่วนนั้นๆ หรือความหนาแน่นของลิงก์ภายในส่วนประกอบนั้น ซึ่งเป็นวิธีการที่ง่ายและไม่ซับซ้อน โดยในงานวิจัยได้กล่าวถึงการวิเคราะห์ในเชิงปริมาณว่าจะพิจารณาในระดับที่สูงขึ้นโดยจะไม่ได้ทำการพิจารณาจากตัวอักษร แต่จะใช้ระดับที่สูงกว่านั้นคือการหาค่าเฉลี่ยความยาวของคำ หรือความยาวของประโยคโดยการพิจารณาด้วยช่องว่าง (White space) หรือมหัพภาค (Full-stop) ในการแบ่งแต่ละคำหรือประโยคออกจากกัน นอกจากนั้นผู้วิจัยได้กล่าวอีกว่าโดยส่วนใหญ่แล้วนั้นเนื้อหาหลักในหน้าเว็บไซต์มักจะถูกล้อมรอบไปด้วย Boilerplate ไม่ว่าจะเป็นส่วนหัวด้านบน (Header) ส่วนท้ายด้านล่าง (Footer) หรือด้านซ้ายและขวาที่มักจะเป็นส่วนของเมนู (Navigation) เพื่อใช้ไปยังหน้าอื่นๆของเว็บไซต์

นอกจากนี้ผู้วิจัยได้ใช้วิธีการคำนวณเพื่อหาความหนาแน่นอีกวิธีหนึ่งที่มีชื่อว่าวิธีการวัดความหนาแน่น (Densitometric Features) ซึ่งเป็นการคำนวณจากการนำจำนวนของตัวอักษรที่อยู่ในบล็อกนั้นๆหารด้วยจำนวนบรรทัด ซึ่งจำนวนบรรทัดจะหาได้จากการนำตัวอักษรในบล็อกนั้นมาตัดทอนด้วยจำนวนตัวอักษรต่อบรรทัดที่เรากำหนด (ในภาษาอังกฤษกำหนดให้อยู่ระหว่าง 80 - 90 ตัวอักษรต่อบรรทัด) สามารถเขียนออกมาเป็นสมการได้ดังนี้

$$T'(b) = \{t | t \in T(l), l_{first}(b) \leq l < l_{last}(b)\}$$

$$g(b) = \begin{cases} \frac{|T'(b)|}{|L(b)|-1} & |L(b)| > 1 \\ |T(b)| & \text{otherwise} \end{cases}$$

โดยที่ $|T(b)|$ คือจำนวนตัวอักษรในบล็อก b หารด้วย $|L(b)|$ คือจำนวนบรรทัดที่หาได้จากการนำตัวอักษรทั้งหมดในบล็อก b มาตัดทอนให้เหลือ 80 - 90 ตัวอักษรต่อบรรทัด ก็จะได้ผลลัพธ์คือความหนาแน่นของบล็อกนั้น ซึ่งจะถูกนำมาพิจารณาเพื่อหาความสัมพันธ์กับการที่จะเป็นเนื้อหาหลักของหน้าเว็บไซต์ต่อไป

งานวิจัยนี้ได้ทำการฝึกฝนตัวจำแนกทั้งสองวิธีที่ได้กล่าวไปคือ การจำแนกจากจำนวนของคำ และการจำแนกจากความหนาแน่นของคำและลิงก์ ซึ่งได้ผลลัพธ์ออกมาเป็นต้นไม้ตัดสินใจที่มีจำนวนโนดเพียง 8 โหนดเท่านั้น สำหรับปัญหาแบบ 2-class (คือการหาเนื้อหา) และ 12 โหนดสำหรับปัญหาแบบ 4-class (คือการหาบทความเต็ม) ทำให้ได้ต้นไม้ตัดสินใจที่เกิดจากการนำชุดข้อมูลฝึกมาพิจารณาจำแนก และทำการหาค่าเกณฑ์แล้วปรับปรุงแก้ไขตัวแปรเหล่านั้น เมื่อเทียบกับอัลกอริทึมที่ใช้การจำแนกแบบคุณลักษณะแบบท้องถิ่น (Local features) คือการจำแนกจากการพิจารณาภายในส่วนประกอบของเว็บไซต์จำพวกเอกสารเว็บไซต์อาจมีความถูกต้องที่น้อยกว่าเล็กน้อย แต่จะเห็นได้ว่าจำนวนของกิ่งของต้นไม้ตัดสินใจนั้นมีจำนวนที่น้อยกว่ามาก มีผลทำให้ความซับซ้อนและความเร็วในการค้นหามีประสิทธิภาพดีขึ้นไปด้วย

Algorithm 1 Densitometric Classifier

```

curr_linkDensity <= 0.333333
| prev_linkDensity <= 0.555556
| | curr_textDensity <= 9
| | | next_textDensity <= 10
| | | | prev_textDensity <= 4: BOILERPLATE
| | | | prev_textDensity > 4: CONTENT
| | | next_textDensity > 10: CONTENT
| | curr_textDensity > 9
| | | next_textDensity = 0: BOILERPLATE
| | | next_textDensity > 0: CONTENT
| prev_linkDensity > 0.555556
| | next_textDensity <= 11: BOILERPLATE
| | next_textDensity > 11: CONTENT
curr_linkDensity > 0.333333: BOILERPLATE

```

ภาพที่ 7 แสดงถึงอัลกอริทึมในแบบจำลองต้นไม้ตัดสินใจของวิธีจำแนกจากความหนาแน่น และความหนาแน่นของลิงก์ [2]

ในตอนทำงานวิจัยได้สรุปว่า ผลการทดลองได้บ่งชี้ว่าข้อความที่มีตัวอักษรที่ยาวกว่า (Long text) มีแนวโน้มที่จะเป็นเนื้อหาหลักมากกว่า เพราะข้อความที่มีตัวอักษรน้อย (Short text) มีแนวโน้มที่จะเป็นส่วนที่ใช้ในการลิงก์ไปยังส่วนอื่นๆของเว็บมากกว่า จึงสามารถสร้างสมมุติฐานได้ว่า เราสามารถจัดส่วนที่มีตัวอักษรหรือประโยคจำนวนน้อยๆ ออกไปได้ เพื่อลดความซับซ้อนและขนาดของข้อมูลก่อนจะนำมาเข้าอัลกอริทึมจำแนก อีกทั้งด้วยตัวอัลกอริทึมเองมีความซับซ้อนที่น้อยและไม่จำเป็นจะต้องสร้างการจำแนกในระดับตัวอักษร (เนื่องจากพิจารณาด้วยจำนวนและความหนาแน่นของคำและประโยคแล้ว) ทำให้มีความเร็วในการจำแนกมากกว่าการใช้อัลกอริทึมที่ยุ่งยากซับซ้อน

2.2.3. Extracting news text from web pages: an application for the visually impaired [5]

งานวิจัยชิ้นนี้ได้กล่าวถึงการนำเทคนิค Boilerplate Detection มาพัฒนาต่อด้วยการนำอัลกอริทึมที่มีชื่อว่าทฤษฎีป่าแบบสุ่ม (Random Forests) เข้ามาช่วยในการพัฒนาเพื่อเพิ่มความแม่นยำให้กับการจำแนกที่แต่ก่อนใช้เพียงต้นไม้ตัดสินใจ (Decision tree) อย่างเดียว ซึ่งยังมีจุดที่สามารถพัฒนาเพื่อเพิ่มความแม่นยำได้อีกด้วยการนำอัลกอริทึมที่มีความซับซ้อนมากขึ้นเพื่อไปเพิ่มศักยภาพให้กับการตัดสินใจ

ผู้วิจัยได้อธิบายถึงการนำทฤษฎีป่าแบบสุ่มมาใช้ในส่วนของ การสร้างตัวตัดสินใจว่าเป็นเนื้อหาหรือไม่ ด้วยการสร้างอาศัยต้นไม้ตัดสินใจจำนวนมากๆ จากการนำชุดข้อมูลไปสุ่มให้มีความแตกต่างกัน เพื่อทำให้เกิดต้นไม้ตัดสินใจที่มีความหลากหลายและนำชุดข้อมูลเข้าไปจำแนกในต้นไม้ตัดสินใจทุกๆต้น ผลลัพธ์ที่ได้ก็คือชุดของต้นไม้ตัดสินใจที่มีประสิทธิภาพ โดยให้ต้นไม้ตัดสินใจแต่ละต้นทำการโหวตว่าข้อมูลที่ได้นำเข้าควรจะถูกจำแนกอยู่ในประเภทใด ในการทดลองของงานวิจัยนี้ได้ทำนำบทความทั้งหมดจำนวน 622 บทความ ซึ่งแต่ละบทความจะประกอบไปด้วย 6 ส่วนคือ

1. หัวข้อ (Heading) คือ ชื่อเรื่องสำคัญของบทความ เป็นสิ่งที่มีอยู่ในทุกๆบทความ
2. ข้อมูลที่เกี่ยวข้องกับบทความ (Supplemental content) คือ ส่วนที่จะอธิบายถึงข้อมูลของบทความ เช่น วันเดือนปีที่แต่ง ผู้แต่ง หรือช่องทางการติดต่อผู้แต่ง เป็นต้น จะมีหรือไม่มีก็ได้
3. ตัวบทความ (Text) คือ ส่วนที่เป็นเนื้อหาสำคัญของบทความ จะต้องมีในทุบบทความ
4. เนื้อหาที่มีเกี่ยวข้อง (Related content) คือ บทความสั้นๆที่อาจจะมีความเกี่ยวข้องกับบทความที่อยู่ในหน้าปัจจุบัน และส่วนนี้มักจะถูกตีความว่าเป็นบทความเช่นเดียวกับตัวบทความ ซึ่งอาจจะมีหรือไม่มีก็ได้
5. ความคิดเห็น (Comments) คือ ส่วนที่ไว้แสดงความคิดเห็นจากผู้อ่านและอาจจะถูกตีความว่าเป็นบทความเช่นเดียวกัน ซึ่งอาจจะมีหรือไม่มีก็ได้
6. ส่วนที่ไม่ต้องการ (Boilerplate) คือ ส่วนอื่นที่ไม่ได้ตกอยู่ในส่วนใดส่วนหนึ่งก่อนหน้า นี้ก็จะถูกตีความว่าเป็นส่วนที่ไม่ต้องการทั้งหมด

ในการเปรียบเทียบจะนำชุดข้อมูลที่ได้จากการแยกส่วนบทความนั้นมาทำการวิเคราะห์และใช้การทดสอบแบบ 10-fold cross validation ผลลัพธ์ที่ได้จะถูกนำไปให้น้ำหนักโดยคิดจากจำนวนคำ

Classifier	Dim	Precision		Recall		F-score		FP rate	
		2-class	4-class	2-class	4-class	2-class	4-class	2-class	4-class
Baseline	0	35.1	32.8	59.2	57.3	44.1	41.7	59.2	57.3
Boilerpipe	6	93.9	89.3	93.8	91.0	93.9	89.5	6.7	8.4
Random forest	6	94.7	92.6	94.7	93.2	94.7	92.9	5.6	6.0
Boilerpipe	64	95.0	91.3	95.0	92.4	95.0	91.4	5.5	7.1
Random Forest	64	96.4	95.1	96.4	95.3	96.4	95.0	3.9	4.4
Random Forest	27	96.6	95.1	96.6	95.3	96.6	95.1	3.5	4.2

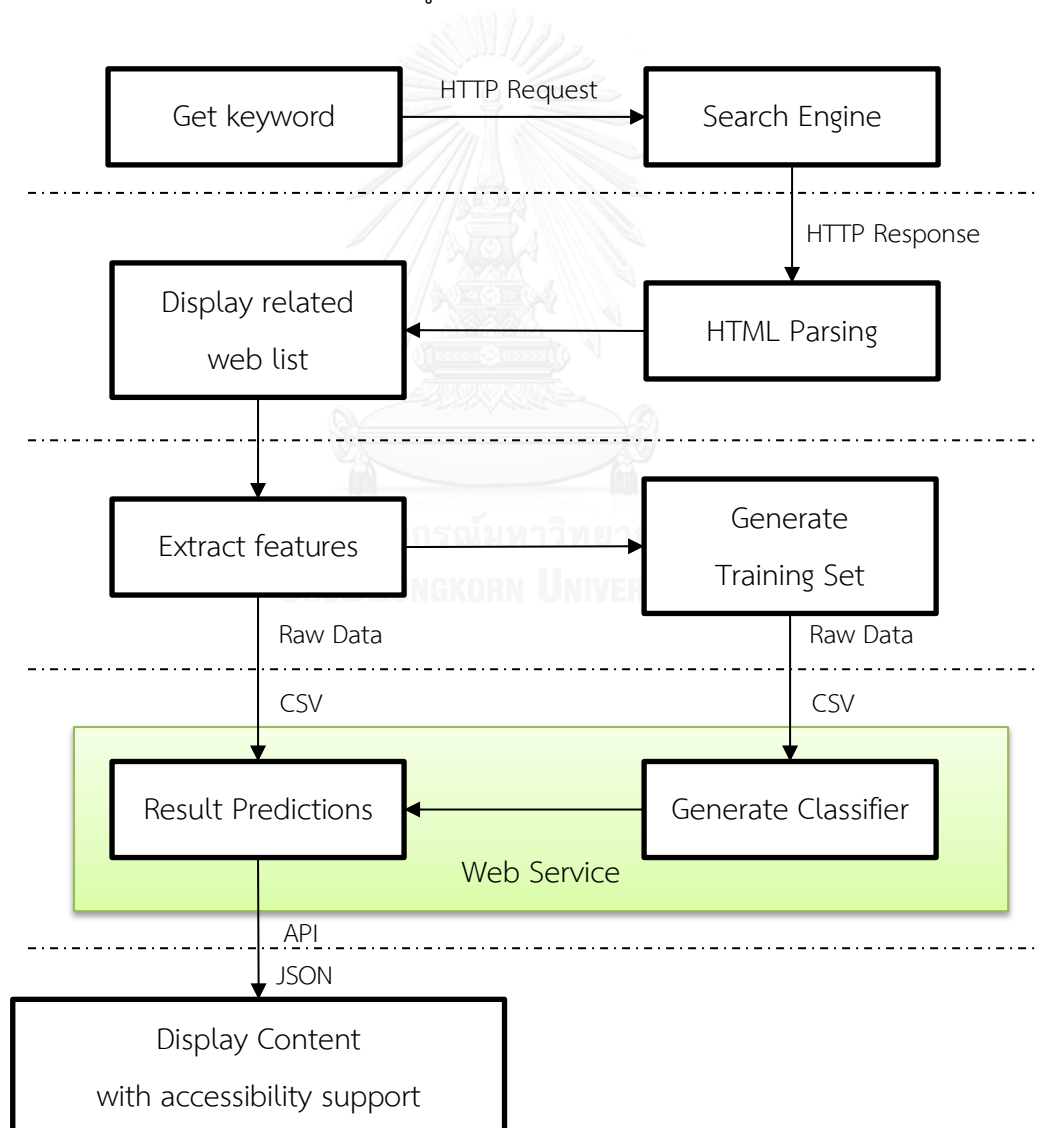
ภาพที่ 8 แสดงการเปรียบเทียบประสิทธิภาพของวิธีการป่าแบบลุ่มกับการใช้ Boilerpipe [5]

ผลลัพธ์ที่ได้จากการทดลองด้วยการเปรียบเทียบกับวิธี Boilerplate Detection จากการใช้ไลบรารีของ Boilerpipe library ได้ผลตามที่แสดงในภาพที่ 8 โดยผู้วิจัยได้สรุปว่าหลังจากที่เราสามารถจำแนกประเภทของข้อความบนหน้าเว็บได้แล้วนั้น เราก็จะสามารถนำไปใช้ประโยชน์ในการจัดรูปแบบที่จะนำเสนอให้กับผู้พิการทางสายตา โดยการนำเอาหัวข้อและเนื้อหาหลักของบทความขึ้นมาอยู่ในตำแหน่งบน เพื่อให้โปรแกรมอ่านหน้าจอสามารถเข้าถึงได้ก่อน แล้วจึงนำส่วนอื่นๆมาใส่ในลำดับต่อมา ถึงแม้ว่าเปอร์เซ็นต์ความถูกต้องที่ได้จะเพิ่มขึ้นมาไม่มาก แต่ก็ถือว่ามีความสำคัญซึ่งทางผู้วิจัยยังได้กล่าวถึงการนำไปใช้ในภาษาอื่นๆ และสามารถยืนยันได้ว่านำไปใช้ได้กับทุกภาษาเนื่องจากตัวอัลกอริทึมของการจำแนกนั้นไม่ได้วิเคราะห์จากภาษา และนับเป็นข้อดีของอัลกอริทึมแบบนี้ ซึ่งใช้วิธีการวิเคราะห์จำนวนคำจากช่องว่างและมหัพภาค จึงทำให้ไม่ขึ้นกับภาษาแต่ก็อาจจะมี ความแม่นยำที่ลดลงเนื่องจากในภาษาอื่นๆอาจจะความแตกต่างในอัตราส่วน ความยาวของ คำ ความยาวของประโยค หรือแม้กระทั่งการใช้มหัพภาคเองก็ตาม โดยแนวทางที่ผู้วิจัยได้บอกถึงคือในส่วนที่เป็น ข้อมูลเกี่ยวกับบทความ (Supplement Content) นั้นยังยากที่จะถูกแยกออกจากตัวเนื้อหาของบทความเนื่องจากไม่มีสิ่งที่สามารถแยกได้จึงเป็นสิ่งที่ควรจะต้องถูกพัฒนาต่อไป

บทที่ 3

แนวคิดและวิธีดำเนินการวิจัย

ในบทนี้จะกล่าวถึงแนวคิดและวิธีการดำเนินการวิจัย โดยแนวคิดและดำเนินการวิจัยในการสร้างเครื่องมือการท่องเว็บไซต์ โดยใช้แบบการค้นหาเนื้อหาหลัก ซึ่งวิธีการดำเนินการวิจัยสามารถแบ่งออกเป็นทั้งหมด 5 ขั้นตอนมีดังต่อไปนี้ 1) การส่งคำร้องขอรายการเว็บไซต์จากคำค้นหา 2) การแปลงข้อมูลเอกสารเว็บไซต์และแสดงผล 3) การสังเคราะห์ข้อมูลเบื้องต้นและการสร้างข้อมูลสำหรับการเรียนรู้จากหน้าเว็บไซต์ 4) การติดต่อเว็บบริการและค้นหาเนื้อหาหลัก และ 5) การแสดงผลสำหรับเครื่องมืออ่านหน้าจอ ซึ่งแสดงอยู่ในแผนภาพดังนี้



ภาพที่ 9 แสดงขั้นตอนการทำงานของเครื่องมือ

3.1. การส่งคำร้องขอรายการเว็บไซต์จากคำค้นหา

เพื่อให้ได้รายการเว็บไซต์ที่มีความเกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการค้นหามากที่สุด ระบบที่ถูกเขียนขึ้นด้วยภาษา Swift จะทำการรับคำค้นหา (Keywords) จากผู้ใช้ ผ่านการพิมพ์ด้วยคีย์บอร์ดที่มีการสนับสนุนจากเครื่องมืออ่านหน้าจอเพื่อทำให้ผู้พิการทางสายตาสามารถพิมพ์ได้อย่างถูกต้อง หรือการใช้โปรแกรมแปลงเสียงเป็นตัวอักษร (Speech-to-text) และนำไปประกอบเป็นคำร้องขอในรูปแบบที่อยู่เว็บไซต์ (URL หรือ Uniform Resource Locator) ร่วมกับการใช้คิวรีสตริง (Query String) และส่งไปยังเว็บไซต์ที่มีระบบค้นหา (Search Engine) โดยโครงสร้างที่อยู่เว็บไซต์เป็นดังนี้

URL: [http://www.google.co.th/search?q=\(W\)&tbm=\(T\)&num=\(N\)&start=\(S\)](http://www.google.co.th/search?q=(W)&tbm=(T)&num=(N)&start=(S))

(W) คือ คำค้นหาที่ได้จากผู้ใช้ โดยเครื่องมือสามารถรองรับคำค้นหาได้ทีละหลายคำ โดยผู้ใช้จะต้องคั่นแต่ละคำด้วยการเว้นวรรค (Whitespace) และเครื่องมือจะทำการแปลงเป็นคิวรีสตริงโดยแทนที่ช่องว่างด้วยเครื่องหมายบวกแทน

(T) คือ ประเภทของการค้นหา ซึ่งประกอบไปด้วยดังนี้ vid หมายถึง การค้นหาวิดีโอที่เกี่ยวข้องกับคำค้นหา isch หมายถึง การค้นหารูปภาพที่เกี่ยวข้องกับคำค้นหา และ nws หมายถึง การค้นหาข่าวที่มีความเกี่ยวข้องกับคำค้นหา โดยในงานวิจัยนี้ได้ใช้การค้นหาประเภท nws

(N) คือ จำนวนการค้นหาต่อการส่งคำขอหนึ่งครั้ง ซึ่งในงานวิจัยนี้ได้กำหนดไว้ที่ 20

(S) คือ ตำแหน่งของการค้นหาเริ่มต้น ซึ่งใช้ในกรณีการค้นหาหน้าถัดไป โดยจะถูกคำนวณจากการนำ (N) ไปคูณกับตำแหน่งหน้าการค้นหา (index) เช่น การค้นหาครั้งแรกกำหนดให้มีค่า (index) เท่ากับ 0 ดังในตัวอย่างที่ 1

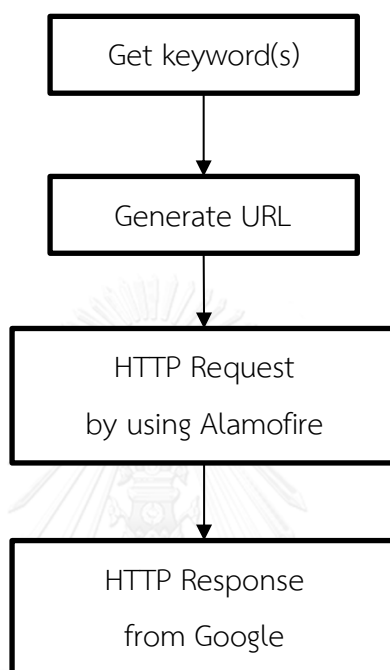
ตัวอย่างที่ 1 ตำแหน่งเริ่มต้น = $(N) \times (\text{index}) = 20 \times 0 = 0$

ตำแหน่งการค้นหาจะมีค่าเท่ากับ 0 ซึ่งมีผลทำให้การแสดงผลรายการเว็บไซต์เริ่มต้น ณ ตำแหน่งที่ 0 ถึง 20 และถ้าหากการค้นหาไปยังหน้าถัดไป ก็จะมีการเพิ่ม (index) ทีละหนึ่ง ซึ่งมีผลเป็นดังตัวอย่างที่ 2

ตัวอย่างที่ 2 ตำแหน่งเริ่มต้น = $(N) \times (\text{index}) = 20 \times 1 = 20$

ทำให้ผลการค้นหาที่ได้มีรายการเว็บไซต์เริ่มต้น ณ ตำแหน่งที่ 20 ถึง 40 เป็นต้น

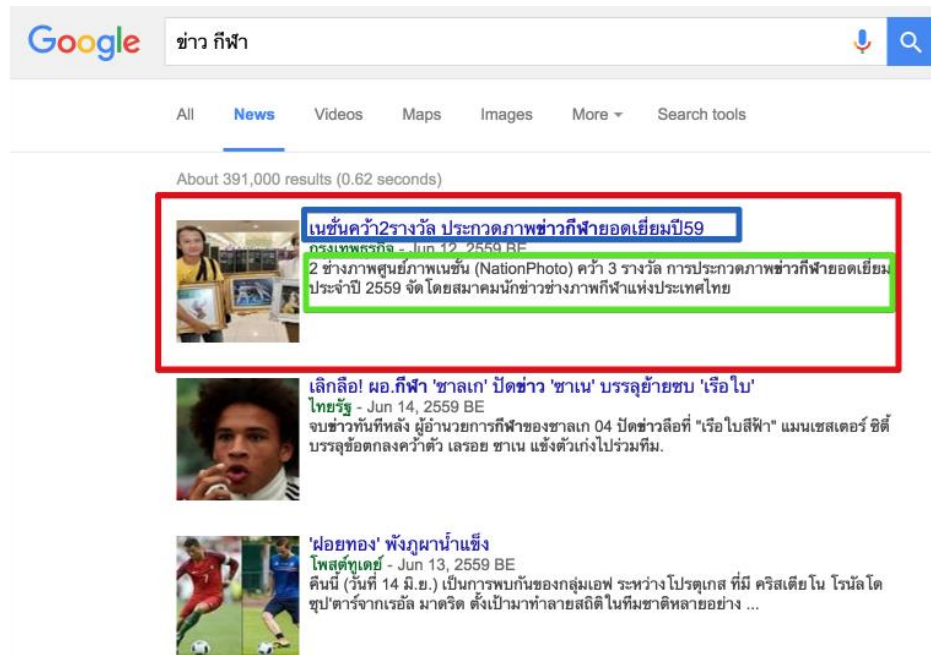
หลังจากเครื่องมือทำการสร้างที่อยู่เว็บไซต์แล้ว ระบบจะทำการส่งรับข้อมูลกับเว็บไซต์ผ่านทางไลบรารีที่มีชื่อว่า อลาโมไฟร์ (Alamofire) หลังจากที่ได้ทำการส่งคำร้องไปยังเว็บไซต์ เครื่องมือจะได้รับข้อมูลกลับมาในรูปแบบเอกสารเว็บไซต์ และเครื่องมือจะทำการแปลงเอกสารเว็บไซต์ โดยแสดงลำดับการทำงานดังในแผนภาพที่ 10



ภาพที่ 10 ขั้นตอนการส่งคำร้องขอรายการเว็บไซต์จากคำค้นหา

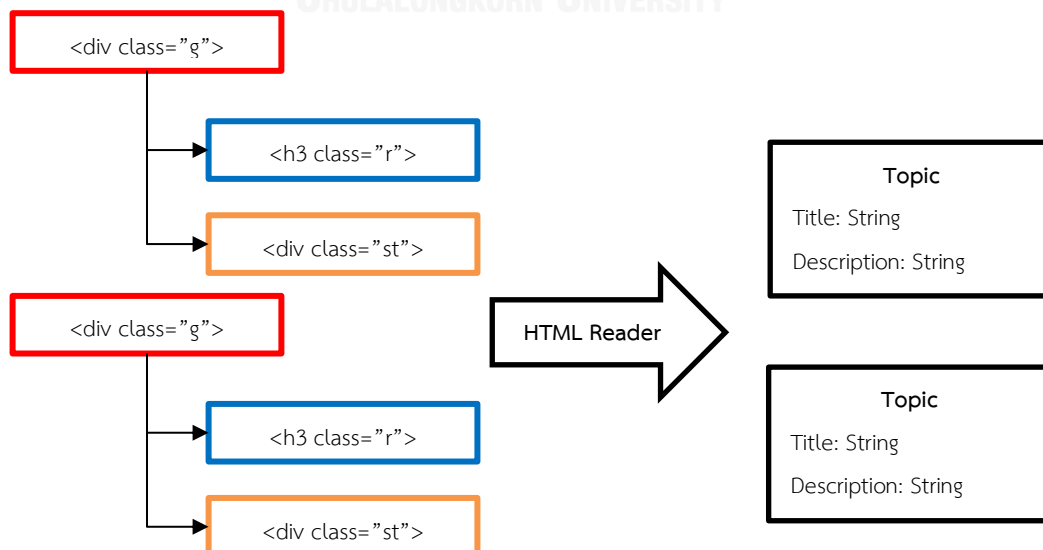
3.2. การแปลงข้อมูลเอกสารเว็บไซต์และแสดงผล

ในการแปลงเอกสารเว็บไซต์ที่ได้มาจากการส่งคำร้องขอในข้อ 3.1 เครื่องมือจะได้รับข้อมูลที่อยู่ในรูปแบบเอกสารเว็บไซต์ในภาษา HTML ดังที่กล่าวในหัวข้อ 2.1.1 จากนั้นเครื่องมือจะใช้ไลบรารีที่ชื่อว่า HTML Reader ทำการแปลงข้อมูลที่ได้ออกมาเป็นเอกสารเว็บไซต์แบบโมเดลต้นไม้ (HTML DOM Tree) ทำให้ระบบสามารถนำข้อมูลที่ถูกรับมาใช้ในการแสดงผลในรูปแบบที่ต้องการได้อย่างง่ายดาย ตรงส่วนนี้ผู้วิจัยได้นำหัวข้อ (Title) และคำอธิบายโดยย่อ (Short Description) ที่ได้จากการค้นหาทำการแสดงผลเป็นรูปแบบของรายการเว็บไซต์ซึ่งถูกจัดเรียงลำดับความเกี่ยวข้องกับคำค้นหาแล้วจากระบบค้นหา พร้อมทั้งพัฒนาให้สามารถสนับสนุนเครื่องมืออ่านหน้าจอด้วยการเพิ่มส่วนที่จะใช้ในการอ่านออกเสียงโดยเครื่องมืออ่านหน้าจอให้สามารถอ่านหัวข้อและคำอธิบายโดยย่อ



ภาพที่ 11 แสดงตัวอย่างผลลัพธ์การค้นหาและนำไปใช้ในการแปลงเอกสาร

จากภาพที่ 11 เป็นการนำคำค้นหาไปทดสอบกับเว็บไซต์กูเกิล ซึ่งในการแปลงเอกสารเว็บไซต์ในส่วนที่เป็นรายการเว็บไซต์ผู้วิจัยได้ทำการคัดแยกส่วนประกอบของรายการเว็บไซต์ โดยในกรอบใหญ่สีแดงคือส่วนที่บรรจุรายละเอียดเว็บไซต์ซึ่งอยู่ในส่วนประกอบ div ที่มีการใช้คลาสที่มีชื่อว่า g เป็นส่วนประกอบ และภายในส่วนประกอบนี้จะมีส่วนสำคัญสองส่วนคือ ส่วนที่อยู่ในกรอบสีน้ำเงินเป็นหัวข้อของแต่ละผลลัพธ์ซึ่งอยู่ในส่วนประกอบ h3 ที่มีคลาสชื่อว่า r และอีกส่วนที่อยู่ในกรอบสีเขียวเป็นรายละเอียดโดยย่อของผลลัพธ์ซึ่งอยู่ในส่วนประกอบ div ที่มีการใช้คลาสที่มีชื่อว่า st ซึ่งสามารถอธิบายในลักษณะรูปภาพดังรูปที่ 12



ภาพที่ 12 แสดงการแปลงเอกสารเว็บไซต์เป็นคลาสที่มีชื่อว่า Topic

3.3. การสังเคราะห์ข้อมูลเบื้องต้นและการสร้างข้อมูลสำหรับการเรียนรู้จากหน้าเว็บไซต์

หลังจากที่ผู้ใช้ทำการเลือกเว็บไซต์จากรายการเว็บไซต์ทั้งหมดแล้ว เครื่องมือจะทำการส่งคำร้องไปยังที่อยู่ที่ถูกเลือก เพื่อที่จะได้เอกสารเว็บไซต์ของหน้าเว็บไซต์นั้นๆ กลับมาและทำการประมวลผลโดยการแปลงข้อมูลเอกสารอีกครั้ง จากนั้นเครื่องมือจะได้ข้อมูลเอกสารเว็บไซต์ ที่อยู่ในรูปแบบโมเดลต้นไม้ และจะถูกนำไปสังเคราะห์ข้อมูลเบื้องต้นโดยใช้หลักการค้นหาคุณสมบัติแบบต้นไม้ ซึ่งได้กล่าวไปแล้วในหัวข้อที่ 2.1.2 จากนั้นทำการเขียนข้อมูลเหล่านี้ลงไฟล์ข้อมูลประเภท Comma Separated Value (CSV) เพื่อที่จะถูกนำไปใช้เป็นข้อมูลสำหรับการเรียนรู้ (Training data set) ของระบบค้นหาเนื้อหาหลัก

ตารางที่ 1 แสดงรายละเอียดข้อมูลที่ถูกสร้างขึ้นจากคุณสมบัติแบบต้นไม้

แอตทริบิวต์	ประเภท
Id	Numeric
p_kw	Numeric
c_kw	Numeric
n_kw	Numeric
p_den	Numeric
c_den	Numeric
n_den	Numeric
p_link	Numeric
c_link	Numeric
n_link	Numeric
Class	Nominal

ในไฟล์ข้อมูลดังกล่าวจะประกอบไปด้วยทั้งหมด 11 คอลัมน์ ซึ่งถูกแบ่งแยกด้วยเครื่องหมายจุลภาค (Comma) และแต่ละแอตทริบิวต์จะถูกคำนวณตามลำดับทิศทางแนวลึก (Depth first) ก่อนที่จะถูกเขียนลงไฟล์ข้อมูล โดยความหมายของแต่ละแอตทริบิวต์จะถูกแสดงอยู่ในตารางที่ 2

ตารางที่ 2 แสดงความหมายของแต่ละแอตทริบิวต์ภายในไฟล์ข้อมูล

แอตทริบิวต์	คำอธิบาย
Id (Identity)	ตัวเลขจำเพาะของแต่ละส่วนประกอบ
p_kw (Previous Element Keyword Count)	จำนวนคำค้นหาที่ปรากฏในส่วนประกอบก่อนหน้าในระดับเดียวกัน
c_kw (Current Element Keyword Count)	จำนวนคำค้นหาที่ปรากฏในส่วนประกอบปัจจุบัน
n_kw (Next Element Keyword Count)	จำนวนคำค้นหาที่ปรากฏในส่วนประกอบถัดไปหน้าในระดับเดียวกัน
p_den (Previous Element Text Density)	ความหนาแน่นของตัวอักษรในส่วนประกอบก่อนหน้าในระดับเดียวกัน
c_den (Current Element Text Density)	ความหนาแน่นของตัวอักษรในส่วนประกอบปัจจุบัน
n_den (Next Element Text Density)	ความหนาแน่นของตัวอักษรในส่วนประกอบถัดไปในระดับเดียวกัน
p_link (Previous Element Link Density)	ความหนาแน่นของตัวอักษรที่เป็นลิงก์ในส่วนประกอบก่อนหน้าในระดับเดียวกัน
c_link (Current Element Link Density)	ความหนาแน่นของตัวอักษรที่เป็นลิงก์ในส่วนประกอบปัจจุบัน
n_link (Next Element Link Density)	ความหนาแน่นของตัวอักษรที่เป็นลิงก์ในส่วนประกอบถัดไปในระดับเดียวกัน
Class	True: เป็นส่วนประกอบที่เป็นเนื้อหาหลัก False: เป็นส่วนประกอบที่ไม่ใช่เนื้อหาหลัก

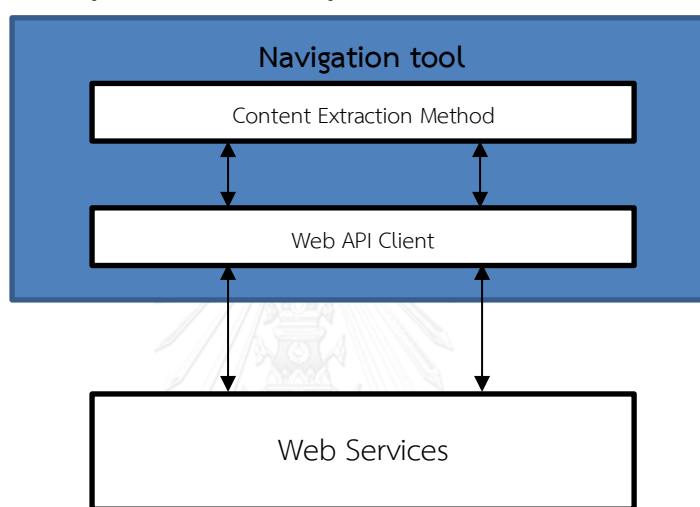
ข่าว ไอที_0_DATA_FEATURE... x										
id	p_kw	c_kw	n_kw	p_den	c_den	n_den	p_link	c_link	n_link	class
1	1	0.0	0.0	8.0	0.0	2.09	15.6	0.00	0.00	0.00, FALSE
2	2	0.0	0.0	0.0	0.0	0.07	1.97	0.00	0.29	0.00, FALSE
3	3	0.0	0.0	0.0	0.0	0.02	0.0	0.00	0.00	0.00, FALSE
4	4	0.0	0.0	0.0	0.0	0.07	1.97	0.0	0.29	0.00, FALSE
5	5	0.0	0.0	0.0	0.0	1.92	0.0	0.00	0.00	0.00, FALSE
6	6	0.0	0.0	0.0	0.0	1.85	0.0	0.00	0.00	0.00, FALSE
7	7	0.0	8.0	0.0	0.0	2.09	15.6	0.0	0.00	0.00, FALSE
8	8	0.0	8.0	0.0	0.0	15.5	0.0	0.00	0.00	0.00, FALSE
9	9	0.0	8.0	0.0	0.0	12.03	2.68	0.00	0.00	0.00, FALSE
10	10	0.0	0.0	3.0	0.0	0.06	6.35	0.00	1.00	0.01, FALSE
11	11	0.0	0.0	0.0	0.0	0.06	0.0	0.00	0.00	0.00, FALSE
12	12	0.0	3.0	3.0	0.06	6.35	3.05	1.00	0.01	0.01, FALSE
13	13	0.0	0.0	3.0	0.0	0.08	6.07	0.00	0.00	0.00, FALSE
14	14	0.0	3.0	0.0	0.08	6.07	0.0	0.00	0.00	0.00, FALSE
15	15	0.0	0.0	0.0	0.0	0.08	0.06	0.00	1.00	1.00, FALSE
16	16	0.0	0.0	0.0	0.0	0.08	0.0	0.00	0.00	0.00, FALSE
17	17	0.0	0.0	0.0	0.0	0.08	0.06	0.08	1.00	1.00, 1.00, FALSE
18	18	0.0	0.0	0.0	0.0	0.06	0.0	0.00	0.00	0.00, FALSE
19	19	0.0	0.0	0.0	0.0	0.06	0.08	0.07	1.00	1.00, 1.00, FALSE
20	20	0.0	0.0	0.0	0.0	0.08	0.0	0.00	0.00	0.00, FALSE
21	21	0.0	0.0	0.0	0.0	0.08	0.07	0.08	1.00	1.00, 1.00, FALSE
22	22	0.0	0.0	0.0	0.0	0.07	0.0	0.00	0.00	0.00, FALSE
23	23	0.0	0.0	3.0	0.07	0.08	4.9	1.00	1.00	0.01, FALSE
24	24	0.0	0.0	0.0	0.0	0.08	0.0	0.00	0.00	0.00, FALSE
25	25	0.0	3.0	0.0	0.08	4.9	0.0	1.00	0.01	0.00, FALSE
26	26	0.0	0.0	3.0	0.0	0.07	4.72	0.00	0.00	0.00, FALSE
27	27	0.0	3.0	0.0	0.07	4.72	0.0	0.00	0.00	0.00, FALSE
28	28	0.0	0.0	0.0	0.0	0.41	0.35	0.00	1.00	1.00, FALSE
29	29	0.0	0.0	0.0	0.0	0.41	0.0	0.00	0.00	0.00, FALSE
30	30	0.0	0.0	0.0	0.0	0.41	0.0	0.00	0.00	0.00, FALSE
31	31	0.0	0.0	0.0	0.0	0.41	0.35	0.4	1.00	1.00, 1.00, FALSE
32	32	0.0	0.0	0.0	0.0	0.35	0.0	0.00	0.00	0.00, FALSE
33	33	0.0	0.0	0.0	0.0	0.35	0.0	0.00	0.00	0.00, FALSE
34	34	0.0	0.0	0.0	0.35	0.4	0.23	1.00	1.00	1.00, FALSE
35	35	0.0	0.0	0.0	0.0	0.4	0.0	0.00	0.00	0.00, FALSE
36	36	0.0	0.0	0.0	0.0	0.4	0.0	0.00	0.00	0.00, FALSE

ภาพที่ 13 แสดงตัวอย่างไฟล์ข้อมูลในรูปแบบ CSV

ภาพที่ 13 คือตัวอย่างไฟล์ข้อมูลที่ถูกสร้างจากการแปลงเอกสารเว็บไซต์เป็นข้อมูลด้วยคุณลักษณะแบบต้นของเว็บไซต์ Siamphone.com (URL: <http://news.siamphone.com/news-27339.html&sa=U&ved=0ahUKEwioofvLycXNAhUM0Y8KHXTYBZYQqQIIIEygAMAA&usg=AFQjCNFA4ST1Cp-PPo5Pv5oycW9Z8vHupA>) ซึ่งเกิดจากการใช้คำค้นหาว่า “ข่าว ไอที” และถูกเลือกจากหัวข้อเว็บไซต์รายการแรกของผลการค้นหา ส่วนที่เป็นคลาส (Class) จะถูกกำหนดให้มีค่าเป็น False มีความหมายว่ามีความเป็นไปได้ว่าจะไม่เป็นเนื้อหาหลัก และเป็น True ถ้ามีความเป็นไปได้ว่าจะ เป็นเนื้อหาหลัก ซึ่งผู้วิจัยจะทำการตรวจสอบจากการนำข้อมูลที่ถูกรับจากการแปลงเอกสารแล้วในอีกไฟล์ข้อมูลหนึ่ง ที่ประกอบไปด้วย id และเนื้อหาภายในส่วนประกอบนั้นๆ จากนั้นนำมาพิจารณาเปรียบเทียบกับไฟล์ข้อมูลแรกด้วยตัวผู้วิจัยเอง ซึ่งหากมีส่วนประกอบที่มีความเป็นไปได้ว่าจะเป็นเนื้อหาหลักก็ทำการเปลี่ยนค่าคลาสของ id นั้นเป็น True ทำให้ได้ไฟล์ข้อมูลที่จะนำไปสร้างไฟล์ข้อมูลเรียนรู้ สุดท้ายไฟล์ข้อมูลนี้จะถูกนำไปเข้าสู่เว็บบริการ เพื่อถูกใช้ในการสร้างตัวจำแนกเนื้อหาต่อไป

3.4. การติดต่อเว็บบริการและค้นหาเนื้อหาหลัก

ในงานวิจัยนี้ได้ใช้เว็บบริการเป็นจุดศูนย์รวมในการวิเคราะห์ข้อมูลที่ได้จากคุณลักษณะแบบต้นในหน้าเว็บไซต์ โดยเครื่องมือจะทำหน้าที่ในการสร้างไฟล์ข้อมูลที่บรรจุข้อมูลในรูปแบบเดียวกับข้อ 3.3 ซึ่งมีลักษณะที่บรรจุข้อมูลดิบในรูปแบบ CSV จากนั้นจะถูกนำไปเข้าระบบของเว็บบริการเพื่อสร้างตัวคัดแยกโดยใช้วิธีการแบบป่าสุ่ม หลังจากการส่งไฟล์ข้อมูลไปที่เว็บบริการแล้วจะต้องมีการจัดการเพื่อแปลงไฟล์ข้อมูลให้เป็นเซตข้อมูล (Data set) ที่สามารถนำไปใช้ในการวิเคราะห์ต่อไปได้ โดยตรงส่วนการจัดการข้อมูลบนเว็บบริการนี้จะถูกอธิบายในหัวข้อ 4.5 ต่อไป



ภาพที่ 14 แสดงรูปแบบการติดต่อเว็บบริการผ่านส่วนต่อประสาน

เครื่องมือจะทำการติดต่อกับเว็บบริการโดยการใช้ไลบรารีจากผู้ให้บริการที่เรียกว่า BigML Swift (ที่อยู่ของไลบรารี <https://github.com/bigmlcom/bigml-swift>) ซึ่งมีการพัฒนาส่วนต่อประสานโปรแกรมประยุกต์ (Application programming interface: API) ที่อยู่ในภาษา Swift อยู่แล้วทำให้ผู้วิจัยสามารถนำส่วนนี้ไปใช้ร่วมในการพัฒนาเครื่องมือที่อยู่ในภาษาเดียวกันได้อย่างง่ายดาย โดยส่วนต่อประสานนั้นมีการพัฒนาให้สามารถรับส่งคำสั่งจากระบบใดๆ ด้วยการส่งคำสั่งแบบไม่ประสานเวลา (Asynchronous) ทำให้ไม่มีการติดขัดของการทำงานของระบบ รวมถึงมีการพัฒนาให้สามารถแปลงผลลัพธ์ที่ได้รับจากการวิเคราะห์ซึ่งมาในรูปแบบ JSON ไปเป็นข้อมูลภายในภาษา Swift ที่สามารถนำออกมาใช้ได้อย่างสะดวกมากขึ้น

ผู้ที่ต้องการพัฒนาเครื่องมือที่จะใช้การประมวลผลการเรียนรู้ของเครื่องจากเว็บบริการ (<http://www.BigML.com>) จะต้องทำการพัฒนาโปรแกรมส่วนต่อประสานร่วมกับการนำข้อมูลที่ต้องการใช้ในการประมวลและทำการอัปโหลดไปสู่เว็บบริการเพื่อให้เว็บบริการทำการประมวลผลแปลงเป็นข้อมูลอยู่ในรูปแบบชุดข้อมูล (Data Set) และจากนั้นจะทำการแปลงเป็นโมเดล (Model)

เพื่อที่จะถูกนำไปสร้างส่วนจำแนกข้อมูล (Classification) ดังรูปภาพที่ 15 จากนั้นจะทำการส่งส่วนจำแนกกลับมาให้เครื่องมือในรูปแบบ JSON ดังรูปภาพที่ 16



ภาพที่ 15 แสดงภาพการทำงานของเว็บบริการ [6]

```

mac$ curl "https://bigml.io/dev/batchprediction/5769764f3bbd213cac006db8?BIGML_AUTH"
{"all_fields": true, "category": 17, "code": 200, "combiner": 1, "confidence": false, "created": "2016-06-21T17:15:59.284000", "credits": 4720.0, "dataset": "dataset/5753c7c93bbd21160c001a58", "dataset_status": true, "description": "", "dev": true, "ensemble": "ensemble/5753c6d13bbd211613000b36", "fields_map": {"000001": "000000", "000002": "000001", "000003": "000002", "000004": "000003", "000005": "000004", "000006": "000005", "000007": "000006", "000008": "000007", "000009": "000008"}, "header": true, "lda": "", "locale": "en-US", "logistic_regression": "", "missing_strategy": 0, "model": "model/5753c6e63bbd21160c001998", "model_status": true, "model_type": 1, "name": "Test_BatchPrediction", "number_of_models": 100, "objective_field": {"column_number": 10, "datatype": "string", "id": "00000a", "name": "class", "optype": "categorical", "order": 8, "preferred": true, "term_analysis": {"enabled": true}}, "output_dataset": true, "output_dataset_resource": "dataset/5769765b3bbd213cb2002af6", "output_dataset_status": true, "output_fields": [], "prediction_name": "result", "private": true, "probabilities": false, "project": "project/56f75ecc3bbd2154910019c8", "resource": "batchprediction/5769764f3bbd213cac006db8?BIGML_AUTH"}
mac$

```

ภาพที่ 16 แสดงผลลัพธ์ในรูปแบบ JSON จากการส่งคำร้องขอไปยังเว็บบริการ

3.5. การแสดงผลสำหรับเครื่องมืออ่านหน้าจอ

การแสดงผลที่ได้จากการค้นหาเนื้อหาหรือการใช้งานต่างๆในเครื่องมือ จะต้องมีการพัฒนาให้มีการรองรับการเข้าถึงด้วยเครื่องมืออ่านหน้าจอ โดยงานวิจัยนี้มีการใช้งานเครื่องมือที่มีชื่อว่า VoiceOver บนอุปกรณ์ระบบไอโอเอส ซึ่งมีหน้าที่ในการอ่านองค์ประกอบต่างๆภายในหน้าจออุปกรณ์ ทำให้ผู้ใช้ที่มีความพิการทางสายตามีความสามารถที่จะได้เข้าถึงเนื้อหาและใช้งานเครื่องมือผ่านทางเสียง

โดยแนวทางการพัฒนาเครื่องมือที่จะทำให้มีการสนับสนุนการเข้าถึง (Accessibility) จากเครื่องมือ VoiceOver ผู้วิจัยจะต้องมีการศึกษาถึงวิธีการออกแบบและพัฒนาโปรแกรมเพิ่มเติมบางส่วนเพื่อที่จะทำให้เครื่องมือของผู้วิจัยมีความสามารถในการอ่านสิ่งต่างๆ ที่ผู้ใช้กำลังสนใจอยู่ ซึ่งมีแอดทริบิวต์ที่ต้องสนใจในการทำให้การแสดงผลรองรับเครื่องมืออ่านหน้าจอ ดังนี้

1. **Label** คือ คำหรือวลีที่อธิบายถึงส่วนควบคุมหรือส่วนแสดงผล แต่ไม่ได้บ่งบอกถึงประเภท เช่น “เพิ่ม” หรือ “เล่น” เป็นต้น
2. **Traits** คือ ส่วนที่ใช้บ่งบอกคุณลักษณะ พฤติกรรม หรือการใช้งานของส่วนประกอบนั้นๆ ตัวอย่างเช่น ส่วนประกอบที่มีพฤติกรรมคล้ายแป้นพิมพ์ ที่ซึ่งมีคุณลักษณะพิเศษจากการรวมแป้นพิมพ์และคุณสมบัติพิเศษ
3. **Hint** คือ วลีโดยย่อซึ่งสามารถอธิบายถึงผลลัพธ์จากการกระทำโดยส่วนประกอบนี้ได้ เช่น “เพิ่มคำนำหน้า” หรือ “เปิดรายการสินค้า” เป็นต้น
4. **Frame** คือ ส่วนที่ใช้บ่งบอกตำแหน่งบนหน้าจอ ซึ่งถูกกำหนดโดยโครงสร้าง CGRect ที่ระบุตำแหน่งของส่วนประกอบจากตำแหน่งบนหน้าจอหรือขนาดได้
5. **Value** คือ ส่วนที่ใช้บ่งบอกคุณสมบัติที่เป็นจำนวนของส่วนประกอบ เมื่อจำนวนนั้นไม่ได้ถูกแสดงผลโดย Label เช่น Label สำหรับตัวเลื่อนที่ใช้ในการปรับค่าชื่อว่า “ความเร็ว” ซึ่งค่าปัจจุบันมีจำนวนอยู่ที่ “50 เปอร์เซ็นต์”

บทที่ 4

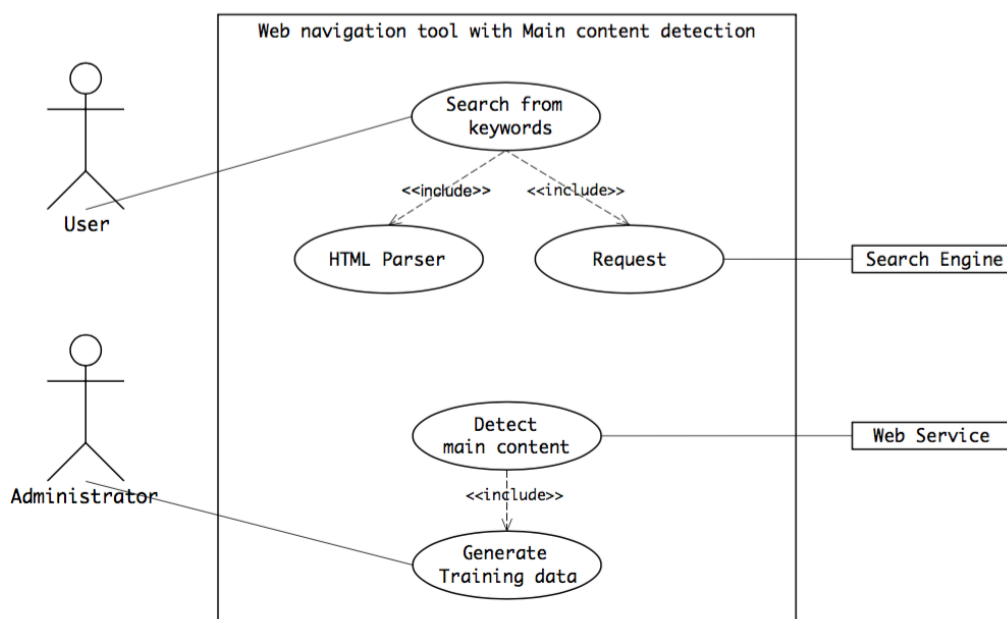
การพัฒนาเครื่องมือ

4.1. ความต้องการเชิงฟังก์ชัน

- 1) ระบบจะต้องสามารถสนับสนุนโปรแกรมอ่านหน้าจอได้อย่างสมบูรณ์
- 2) ระบบจะต้องสามารถรับคำค้นหาจากผู้ใช้ได้ด้วยการใช้คีย์บอร์ดหรือการใช้การแปงเสียงเป็นตัวอักษรได้
- 3) ระบบจะต้องสามารถแปลงเอกสารเว็บไซต์ที่ได้จากการค้นหาให้อยู่ในรูปแบบที่ถูกต้องบนเครื่องมือได้
- 4) ระบบจะต้องสามารถติดต่อกับเว็บบริการเพื่อใช้ค้นหาเนื้อหาหลักได้
- 5) ระบบจะต้องสามารถให้ผู้ดูแลระบบสามารถพัฒนาข้อมูลเรียนรู้บนเว็บบริการได้
- 6) ระบบจะต้องสามารถค้นหาเนื้อหาหลักจากหน้าเว็บไซต์ที่ถูกเลือกโดยผู้ใช้ได้อย่างถูกต้อง

4.2. การวิเคราะห์ความต้องการและแผนภาพฟังก์ชันงานของระบบ

ซึ่งจากหัวข้อ 4.1 เราสามารถนำมาพิจารณาและวิเคราะห์ความต้องการเชิงฟังก์ชันที่ได้ เพื่อนำมาออกแบบเครื่องมือซึ่งจะมีการทำงานของระบบเป็นไปตามที่แสดงเป็นแผนภาพยูสเคส (Use case diagram) ดังในภาพที่ 17



ภาพที่ 17 แผนภาพยูสเคสแสดงความสัมพันธ์ของระบบ

จากแผนภาพยูสเคสดังกล่าว แสดงให้เห็นถึงความเกี่ยวข้องของผู้ใช้งานเครื่องมือ การตอบสนองของเครื่องมือ และระบบการจัดการที่ใช้สร้างตัวค้นหาเนื้อหาหลักบนเว็บบริการที่ใช้ในงานวิจัยนี้ โดยมีรายละเอียดของแต่ละฟังก์ชันในการใช้งานดังต่อไปนี้

1) ค้นหาเว็บไซต์จากคำค้นหา

ผู้ใช้งานทำการป้อนคำค้นหาให้กับเครื่องมือ โดยผ่านทางกรใช้เครื่องมือป้อนคำเช่น แป้นพิมพ์ที่มีการสนับสนุนเครื่องมืออ่านหน้าจอ ที่จะทำหน้าที่อ่านตัวอักษรที่อยู่บนแป้นขณะนั้นออกมาเป็นเสียง หรือการใช้เครื่องมือแปลงเสียงเป็นตัวอักษร ซึ่งคำค้นหาที่ได้จะถูกนำไปแปลงเป็นที่อยู่เว็บไซต์และส่งคำร้องขอไปยังผู้ให้บริการ ซึ่งในงานวิจัยนี้ผู้ให้บริการคือเว็บไซต์กูเกิลเพื่อให้ได้เอกสารเว็บไซต์และนำมาแปลงให้อยู่ในรูปแบบที่ต้องการ

2) แปลงเอกสารเว็บไซต์

ในเครื่องมือจะมีระบบการแปลงเอกสารเว็บไซต์สำหรับการแปลงเอกสารที่อยู่ในภาษา HTML มาเป็นรูปของอ็อบเจกต์ใหม่เพื่อนำไปทำการประมวลผล และสร้างวิธีการนำเสนอให้กับผู้ใช้งานในรูปแบบใหม่ที่มีการรองรับจากเครื่องมืออ่านหน้าจอ หรือแม้กระทั่งการแปลงเอกสารเว็บไซต์เป็นไฟล์ข้อมูลจากการคำนวณส่วนประกอบแต่ละส่วนเพื่อหาค่าคุณสมบัติแบบต้น แล้วจึงทำการเขียนข้อมูลเหล่านั้นลงเป็นไฟล์ข้อมูลชนิด CSV เพื่อที่จะนำไปประมวลผลในเว็บบริการต่อไป

3) สร้างข้อมูลสำหรับการเรียนรู้

หลังจากการแปลงเอกสารเว็บไซต์แล้วทำการสร้างไฟล์ข้อมูลที่ได้จากการคำนวณจากคุณสมบัติแบบต้นมาแล้ว ผู้ดูแลระบบจะทำการสร้างข้อมูลให้กับไฟล์เหล่านั้นจากการเปรียบเทียบกับข้อมูลจริง และทำการอัปโหลดไฟล์ดังกล่าวไปเว็บบริการเพื่อทำการประมวลผล และสร้างชุดข้อมูลที่จะนำไปใช้ในการสร้างวิธีการป่าแบบสุ่ม (Random forest) ซึ่งเมื่อสามารถสร้างป่าแบบสุ่มออกมาได้แล้ว เว็บบริการก็จะมีหน้าที่ในการวิเคราะห์ข้อมูลเว็บไซต์ที่ถูกเลือกจากผู้ใช้ และค้นหาส่วนประกอบที่มีความเป็นไปได้ที่จะเป็นเนื้อหาหลักส่งกลับมายังเครื่องมือ เพื่อทำการแสดงผลพร้อมรองรับให้เครื่องมืออ่านหน้าจอทำการอ่านเนื้อหาที่ได้ออกมาเป็นเสียงแก่ผู้ใช้ทันที

4.3. สภาพแวดล้อมที่ใช้ในการพัฒนาเครื่องมือสนับสนุน

สภาพแวดล้อมที่ใช้ในการพัฒนาระบบจะอ้างอิงมาจากคอมพิวเตอร์ที่ใช้ในการพัฒนา โดยประกอบไปด้วยฮาร์ดแวร์ (Hardware) และซอฟต์แวร์ (Software) ที่ใช้ในการพัฒนาระบบ ซึ่งมีรายละเอียดดังนี้

1) ระบบฮาร์ดแวร์

เครื่องคอมพิวเตอร์ที่ใช้ในการพัฒนาระบบควรมีฮาร์ดแวร์ขั้นต่ำดังต่อไปนี้

- หน่วยการประมวลผล (CPU) Core i5 ความเร็ว 2.6 กิกะเฮิร์ตซ์ (intel(R) Core(TM) i5 CPU 2.6 GHZ)
- หน่วยความจำสำรอง (Memory) ความเร็ว 8 กิกะไบต์ (Ram 8 GB)
- งานบันทึกแบบแข็ง (Hard disk) ความจุ 128 กิกะไบต์ (Hard disk 128 GB)

2) ซอฟต์แวร์

เครื่องคอมพิวเตอร์ที่ใช้ในการพัฒนาระบบมีซอฟต์แวร์ดังต่อไปนี้

- ระบบปฏิบัติการโอเอสเอ็กซ์ เวอร์ชัน 10 (OSX El Capitan 10.11.5)
- โปรแกรม XCode เวอร์ชัน 7
- กูเกิลโครม เว็บเบราว์เซอร์ (Google Chrome web browser) เวอร์ชัน 46.0.2490.86 m

3) การติดตั้งซอฟต์แวร์

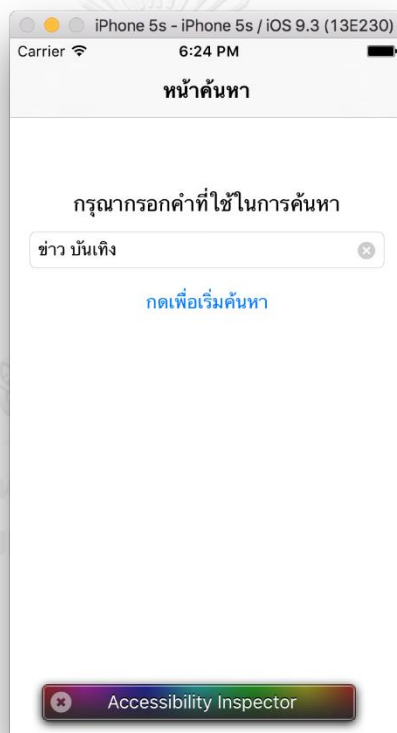
ทำการติดตั้งเครื่องมือในการพัฒนาระบบทั้งหมดลงในเครื่องคอมพิวเตอร์ที่ใช้พัฒนาระบบโดยเริ่มลำดับการติดตั้งตามขั้นตอนดังต่อไปนี้

- 1) ติดตั้งระบบปฏิบัติการโอเอสเอ็กซ์ (OSX El Capitan)
- 2) ติดตั้งโปรแกรม XCode เวอร์ชัน 7 ขึ้นไป
- 3) ติดตั้งโปรแกรมส่วนเสริม Cocoapod
- 4) ติดตั้งโปรแกรมส่วนเสริม Alamofire
- 5) ติดตั้งโปรแกรมส่วนเสริม HTML Reader
- 6) ติดตั้งโปรแกรม Google Chrome

4.4. ขั้นตอนการทำงานของเครื่องมือ

4.4.1. หน้าค้นหา

ในหน้าค้นหา จะประกอบไปด้วยกล่องข้อความที่จะใช้ในการป้อนคำค้นหาจากผู้ใช้ และปุ่มค้นหา ซึ่งในกล่องข้อความนี้ผู้ใช้สามารถป้อนได้มากกว่าหนึ่งคำ โดยการใช้การเว้นวรรคในการแยกคำค้นหาออกจากกัน โดยผู้วิจัยได้กำหนดให้คำค้นหาแรกมีความสำคัญมากที่สุด และเรียงลำดับลงมาตามจำนวนคำค้นหา ความสำคัญนี้ไม่ได้มีผลต่อการค้นหารายการเว็บไซต์ที่เกี่ยวข้อง แต่จะถูกนำมาใช้เป็นส่วนหนึ่งในคุณสมบัติแบบต้นที่ผู้วิจัยต้องการทดสอบแทน ซึ่งในหน้าค้นหาจะมีการออกแบบให้ดูเรียบง่ายเพื่อลดความซับซ้อนในการป้อนคำค้นหาจากผู้ใช้ที่ไม่สามารถมองเห็น และอาศัยเครื่องมืออ่านหน้าจอในการใช้งาน ดังที่แสดงในภาพที่ 18



ภาพที่ 18 แสดงหน้าจอเครื่องมือในส่วนของหน้าค้นหา

4.4.2. หน้าผลลัพธ์การค้นหา (รายการเว็บไซต์)

ซึ่งหลังจากผู้ใช้กดปุ่มค้นหาแล้วเครื่องมือจะทำการสร้าง URL ที่ประกอบไปด้วยคำค้นหา และประเภทการค้นหาซึ่งได้กล่าวไว้ในหัวข้อ 3.1 ผ่านทางการใช้ส่วนเสริมที่มีชื่อว่า อลาโมฟารีย์ในการส่งคำร้องขอไปยังที่อยู่เว็บไซต์และทำการรับผลลัพธ์การร้องขอที่อยู่ในรูปแบบเอกสารเว็บไซต์ประเภทข้อความตัวอักษร

จากนั้นเครื่องมือจะทำการแปลงเอกสารเว็บไซต์ให้อยู่ในรูปแบบของอ็อบเจกต์โมเดลต้นไม้ และเก็บไว้ในตัวแปร ทำให้เครื่องมือสามารถนำไปใช้ในการประมวลผลเพื่อใช้ในการสร้างรายการเว็บไซต์ โดยจะทำการค้นหาส่วนประกอบที่เป็นหัวข้อและเนื้อหาโดยย่อจากเอกสารเว็บไซต์และแปลงให้อยู่ในรูปแบบโมเดลของหัวข้อ (Topic) ที่ประกอบไปด้วยหัวข้อเรื่อง (Title) และเนื้อหาโดยย่อ (Description) ซึ่งในการแสดงผลผู้วิจัยได้ใช้หัวข้อเรื่องในการแสดงในรายการเว็บไซต์ ดังภาพที่ 19 และเมื่อผู้ใช้ใช้เครื่องมืออ่านหน้าจอไปสัมผัส ฅ รายการใดรายการหนึ่งเครื่องมืออ่านหน้าจอจะอ่านหัวข้อเรื่องและเนื้อหาโดยย่อขึ้นมา เพื่อให้ผู้ใช้สามารถเข้าถึงเนื้อหาโดยย่อได้ว่ามีความเกี่ยวข้องกับสิ่งที่ผู้ใช้ต้องการหรือไม่ก่อนที่จะเข้าถึงหน้าเนื้อหาหลักแบบสมบูรณ์



ภาพที่ 19 แสดงหน้าจอเครื่องมือในส่วนของหน้าผลลัพธ์การค้นหา

4.4.3. การคัดแยกเนื้อหา

หลังจากที่ผู้ใช้ได้ทำการเลือกเว็บไซต์จากรายการเว็บไซต์แล้ว เครื่องมือจะทำการส่งที่อยู่ของเว็บไซต์เพื่อทำการส่งคำร้องขอไปยังที่อยู่เว็บไซต์นั้นและรับเอกสารเว็บไซต์อีกครั้ง จากนั้นทำการแปลงเอกสารเว็บไซต์ให้อยู่ในรูปแบบของอ็อบเจกต์ แต่จากจุดนี้เครื่องมือจะทำการคำนวณที่ละส่วนประกอบบนหน้าเว็บไซต์เพื่อให้ได้มาซึ่งคุณลักษณะแบบต้นโดยเป็นไปตามลำดับแบบลึกก่อน (Depth first traversal) ที่เกิดจากการเขียนโปรแกรมแบบเรียกซ้ำ (Recursive) โดยการคำนวณจะมี

การสร้างตัวเลขที่มีลักษณะจำเพาะหรือที่เรียกว่าไอดี (ID: identity) ของแต่ละส่วนประกอบเพื่อที่จะนำไปใช้ในการสร้างทั้งไฟล์ข้อมูลสำหรับการเรียนรู้ และไฟล์ข้อมูลจริงเพื่อใช้ในการส่งไปให้เว็บบริการเพื่อค้นหาเนื้อหาหลัก ซึ่งไฟล์ข้อมูลสำหรับการเรียนรู้ (Training Set) จะมีการค้นหาเนื้อหาด้วยตัวผู้วิจัยเอง เพื่อใช้ในการกำหนดว่าส่วนประกอบใดที่จะเป็นเนื้อหาหลักของหน้าเว็บไซต์ โดยการกำหนดให้ค่าของคลาสมีค่าเป็น TRUE และในทางกลับกันถ้าเป็นส่วนที่ไม่ใช่เนื้อหา ค่าของคลาสจะถูกกำหนดให้เป็น FALSE ดังในรูป

id	p_kw	c_kw	n_kw	p_den	c_den	n_den	p_link	c_link	n_link	class
1	0	0	100	0	0.56	33.5	0	0	0	FALSE
2	0	0	0	0	0.49	0	0	0	0	FALSE
3	0	0	0	0	0.1	0.11	0	0	0	FALSE
4	0	0	0	0.1	0.11	0	0	0	0	FALSE
5	0	100	42	0.56	33.5	161.57	0	0	0	FALSE
6	0	100	0	0	33.43	0	0	0	0	FALSE
7	0	50	50	0	16.61	16.68	0	0	0	FALSE
8	0	50	0	0	16.51	0	0	0	0	FALSE
9	0	50	0	0	16.51	0	0	0	0	FALSE
10	0	0	6	0	0.06	0.36	0	1	0.19	FALSE
11	0	0	0	0	0.06	0	0	0	0	FALSE
12	0	6	10	0.06	0.36	0.85	1	0.19	0.09	FALSE
13	0	1	5	0	0.07	0.27	0	0	0	FALSE
14	1	5	0	0.07	0.27	0	0	0	0	FALSE
15	0	3	2	0	0.1	0.12	0	1	1	FALSE
16	0	3	0	0	0.1	0	0	0	0	FALSE
17	3	2	0	0.1	0.12	0	1	1	0	FALSE
18	0	2	0	0	0.12	0	0	0	0	FALSE
19	6	10	8	0.36	0.85	0.64	0.19	0.09	0.11	FALSE
20	0	0	10	0	0.08	0.75	0	0	0	FALSE
21	0	10	0	0.08	0.75	0	0	0	0	FALSE
22	0	2	2	0	0.1	0.19	0	1	1	FALSE
23	0	2	0	0	0.1	0	0	0	0	FALSE
24	2	2	2	0.1	0.19	0.11	1	1	1	FALSE
25	0	2	0	0	0.19	0	0	0	0	FALSE
26	2	2	2	0.19	0.11	0.15	1	1	1	FALSE
27	0	2	0	0	0.11	0	0	0	0	FALSE
28	2	2	2	0.11	0.15	0.09	1	1	1	FALSE
29	0	2	0	0	0.15	0	0	0	0	FALSE
30	2	2	0	0.15	0.09	0	1	1	0	FALSE

ภาพที่ 20 ตัวอย่างไฟล์ข้อมูลในรูปแบบ CSV ที่มีการระบุคลาส

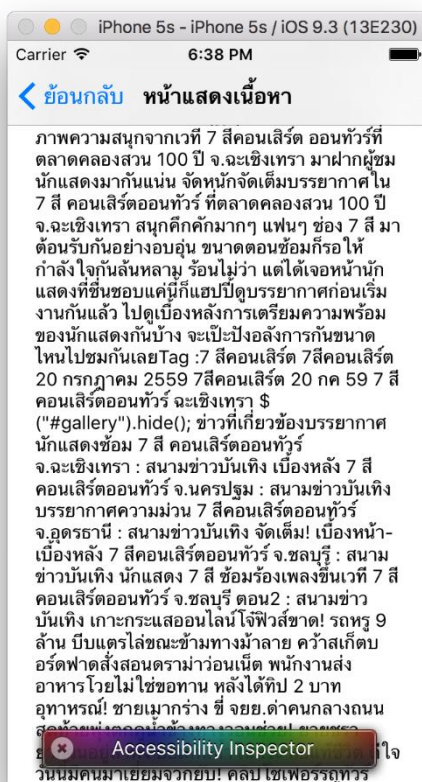
ซึ่งไฟล์ข้อมูลเหล่านี้จะถูกสร้างจากที่อยู่เว็บไซต์ที่เลือกและนำมาประกอบกันเป็นไฟล์สำหรับการเรียนรู้เพื่อใช้ในการสร้างตัวจำแนกเนื้อหาจากวิธีการป่าแบบสุ่ม ซึ่งส่วนนี้จะถูกสร้างโดยผู้วิจัยก่อนที่จะถูกนำไปใช้ในการสร้างตัวจำแนก โดยไฟล์ข้อมูลจริงนั้นจะใช้วิธีการเดียวกันกับไฟล์ข้อมูลเรียนรู้ แต่จะไม่มีกระบวนการระบุแอตทริบิวต์คลาส เนื่องจากแอตทริบิวต์นี้ถูกสร้างจากการนำไฟล์ข้อมูลจริงผ่านการคัดแยกเนื้อหาด้วยวิธีการป่าแบบสุ่มที่สร้างไว้ และสุดท้ายจะได้ชุดข้อมูลที่ผ่านกระบวนการจำแนก (Batch Prediction) ซึ่งมีการระบุคลาสว่าเป็นเนื้อหาหลักหรือไม่อยู่บนเว็บบริการ ในขั้นตอนสุดท้ายเครื่องมือจะทำการส่งคำร้องขอเพื่อให้เว็บบริการส่งคำตอบกลับที่อยู่ในรูปแบบไฟล์ข้อมูล CSV กลับมาให้ จากนั้นก็จะนำไฟล์ข้อมูลที่ได้ไปทำการประมวลผลในขั้นต่อไป

id	p_kw	c_kw	n_kw	p_den	c_den	n_den	p_link	c_link	n_link	class	result	
1	0	0	100	0	0.56	33.5	0	0	0	0	FALSE	FALSE
2	0	0	0	0	0.49	0	0	0	0	0	FALSE	FALSE
3	0	0	0	0	0.1	0.11	0	0	0	0	FALSE	FALSE
4	0	0	0	0.1	0.11	0	0	0	0	0	FALSE	FALSE
5	0	100	42	0.56	33.5	161.57	0	0	0	0	FALSE	FALSE
6	0	100	0	0	33.43	0	0	0	0	0	FALSE	FALSE
7	0	50	50	0	16.61	16.68	0	0	0	0	FALSE	FALSE
8	0	50	0	0	16.51	0	0	0	0	0	FALSE	FALSE
9	0	50	0	0	16.51	0	0	0	0	0	FALSE	FALSE
10	0	0	6	0	0.06	0.36	0	1	0.19	0	FALSE	FALSE
11	0	0	0	0	0.06	0	0	0	0	0	FALSE	FALSE
12	0	6	10	0.06	0.36	0.85	1	0.19	0.09	0	FALSE	FALSE
13	0	1	5	0	0.07	0.27	0	0	0	0	FALSE	FALSE
14	1	5	0	0.07	0.27	0	0	0	0	0	FALSE	FALSE
15	0	3	2	0	0.1	0.12	0	1	1	1	FALSE	FALSE
16	0	3	0	0	0.1	0	0	0	0	0	FALSE	FALSE
17	3	2	0	0.1	0.12	0	1	1	1	0	FALSE	FALSE
18	0	2	0	0	0.12	0	0	0	0	0	FALSE	FALSE
19	6	10	8	0.36	0.85	0.64	0.19	0.09	0.11	0	FALSE	FALSE
20	0	0	10	0	0.08	0.75	0	0	0	0	FALSE	FALSE
21	0	10	0	0.08	0.75	0	0	0	0	0	FALSE	FALSE
22	0	2	2	0	0.1	0.19	0	1	1	1	FALSE	FALSE
23	0	2	0	0	0.1	0	0	0	0	0	FALSE	FALSE
24	2	2	2	0.1	0.19	0.11	1	1	1	1	FALSE	FALSE
25	0	2	0	0	0.19	0	0	0	0	0	FALSE	FALSE
26	2	2	2	0.19	0.11	0.15	1	1	1	1	FALSE	FALSE
27	0	2	0	0	0.11	0	0	0	0	0	FALSE	FALSE
28	2	2	2	0.11	0.15	0.09	1	1	1	1	FALSE	FALSE
29	0	2	0	0	0.15	0	0	0	0	0	FALSE	FALSE
30	2	2	0	0.15	0.09	0	1	1	0	0	FALSE	FALSE

ภาพที่ 21 ตัวอย่างผลลัพธ์การคาดเดาข้อมูลที่ได้จากเว็บบริการ

4.4.4 หน้าแสดงผลเนื้อหา

หลังจากที่เครื่องมือส่งคำร้องไปยังเว็บบริการเพื่อให้ได้ผลลัพธ์ดังในรูปภาพที่ 21 แล้วจากนั้น จะทำการเปรียบเทียบไอดีระหว่างชุดข้อมูล (Data set) กับข้อมูลที่ได้มาจากเอกสารเว็บไซต์ เพื่อนำ ส่วนประกอบที่มีไอดีตรงกันมาแสดงผลเป็นเนื้อหาหลักบนหน้าจอแสดงผลดังที่แสดงในภาพที่ 22



ภาพที่ 22 หน้าแสดงเนื้อหาหลักบนเครื่องมือ

4.5 ขั้นตอนการจัดการระบบเว็บบริการ

ในงานวิจัยนี้ได้ใช้ระบบเว็บบริการในการสร้างตัวคัดแยก (Classifier) โดยใช้ข้อมูลจากเครื่องมือที่สร้างขึ้นผนวกกับการใช้ส่วนต่อประสานของผู้พัฒนาเว็บบริการ ซึ่งทำให้การทำงานของเครื่องมือในการคัดแยกเนื้อหาสามารถเป็นไปอย่างง่ายตายและลดการใช้ทรัพยากรบนเครื่องมืออุปกรณ์ของผู้ใช้เอง แต่ก่อนที่จะสามารถใช้งานเว็บบริการได้จะต้องมีการจัดการข้อมูลเบื้องต้นเพื่อทำการสร้างตัวคัดแยกจากวิธีการป่าแบบสุ่มขึ้นมาก่อน ด้วยการอาศัยไฟล์ข้อมูลที่ถูกสร้างจากผู้วิจัยในหัวข้อ 4.4.3 ซึ่งเป็นไฟล์ข้อมูลเรียนรู้ที่อยู่ในรูปแบบ CSV

The screenshot shows the 'Sources' page in the bigml interface for project 'ATapC'. The page displays a table of sources with columns for Type, Name, Last Updated, Size, and a status icon. A tooltip is visible over the 'training_set.csv' row, stating 'Create a source from a local file (.csv, .tsv, .txt, .json, .arff, .data, .gz, .bz2)'. The table contains the following data:

Type	Name	Last Updated	Size	Status
CSV	test_set.csv	4d 2h	99.0 KB	1
CSV	training_set.csv	4d 3h	1.1 MB	1
CSV	training_set_new.csv	4d 3h	391.3 KB	1
CSV	testset_noid.csv	1m	141.0 KB	1
CSV	trainingset.csv	1m 2w	509.3 KB	1

ภาพที่ 23 แสดงหน้าเว็บบริการในส่วนของการสร้างไฟล์เรียนรู้

จากรูปภาพที่ 23 จะพบว่าหลังจากที่ผู้วิจัยทำการส่งไฟล์ข้อมูลเรียนรู้ไปให้ทางเว็บบริการแล้ว จะต้องมีส่วนที่แปลงไฟล์ข้อมูลดังกล่าวให้อยู่ในรูปแบบที่เว็บบริการสามารถนำไปใช้ในการสร้างตัวคัดแยกด้วย

The screenshot shows the 'Sources' page in the bigml interface for project 'ATapC'. A context menu is open over the 'training_set.csv' source, displaying the following options:

- 1-CLICK DATASET
- VIEW DETAILS
- DELETE SOURCE
- MOVE TO...

ภาพที่ 24 แสดงหน้าเว็บบริการในส่วนของการสร้างชุดข้อมูลจากไฟล์ข้อมูล

หลังจากที่ได้ทำการอัปโหลดไฟล์ข้อมูลเรียนรู้แล้ว ผู้วิจัยจะต้องทำการสร้างชุดข้อมูล (Data Set) จากการนำไฟล์ข้อมูลเรียนรู้ดังกล่าวไปสร้างเป็นชุดข้อมูลใหม่ ดังที่แสดงในรูปภาพที่ 24

The screenshot displays the bigml web interface for a project named 'ATapC'. The top navigation bar includes 'PRIVATE DEPLOYMENTS', 'GALLERY', 'LABS', 'GMKIPLADEN', 'DOCUMENTATION', 'HELP & SUPPORT', and 'Dashboard'. The main navigation bar shows 'Sources', 'Datasets', 'Models', 'Clusters', 'Anomalies', 'Associations', 'Predictions', and 'Tasks'. The 'Datasets' tab is active, showing a table of fields for the 'trainingset dataset'. A context menu is open over the table, listing various configuration options.

Name	Type	Count	Missi
id	123	14,786	0
p_kw	123	14,786	0
c_kw	123	14,786	0
n_kw	123	14,786	0
p_den	123	14,786	0
c_den	123	14,786	0
n_den	123	14,786	0
p_link	123	14,786	0
c_link	123	14,786	0
n_link	123	14,786	0

The context menu options are:

- CONFIGURE MODEL
- CONFIGURE ENSEMBLE
- CONFIGURE LOGISTIC REGRESSION (NEW)
- CONFIGURE CLUSTER
- CONFIGURE ANOMALY
- CONFIGURE ASSOCIATION
- TRAINING AND TEST SET SPLIT
- SAMPLE DATASET
- FILTER DATASET
- ADD FIELDS TO DATASET

ภาพที่ 25 แสดงหน้าเว็บบริการในส่วนของการสร้างตัวจำแนกด้วยป่าแบบสุ่ม

ภาพที่ 26 แสดงหน้าเว็บบริการในส่วนของการกำหนดคุณลักษณะของการสร้างตัวจำแนก

จากนั้นจะต้องทำการสร้างตัวคัดแยกจากชุดข้อมูลที่ได้มาก่อนหน้านี้ ซึ่งเราสามารถกำหนดขนาดจำนวนต้นไม้ตัดสินใจ การเพิ่มกระบวนการในการสร้างสมดุลให้กับคลาส (Class Balancer) หรือกำหนดน้ำหนักของแต่ละแอตทริบิวต์ได้ (Objective Weight) โดยในงานวิจัยนี้ได้กำหนดให้สร้างต้นไม้ตัดสินใจแบบสุ่มทั้งหมด 128 ต้น ร่วมกับการสร้างสมดุลให้กับคลาส เนื่องจากมาจากจำนวนคลาสที่เป็น TRUE มีน้อยกว่ามากจึงจำเป็นต้องมีการกำหนดให้มีการสร้างค่าสมดุล และกำหนดให้ใช้วิธีแบบการสุ่มต้นไม้ตัดสินใจ (Random Decision Forest) ใช้ชุดข้อมูลทดสอบจากการนำเว็บไซต์ทั้งหมด 54 เว็บไซต์จากคำค้นหาที่แตกต่างกันทั้งสิ้น 10 คำมีข้อมูลส่วนประกอบจำนวนทั้งหมด 23,417 ตัว มีคลาสที่เป็นเนื้อหาทั้งสิ้น 56 ตัวและ 23,361 ตัวที่เป็นส่วนประกอบอื่นๆ

สุดท้ายเมื่อได้ตัวคัดแยกแล้ว ผู้วิจัยจะสามารถทำการทดสอบตัวคัดแยกดังกล่าวได้โดยการสร้างชุดไฟล์ข้อมูลทดสอบ (Test set) ที่สร้างโดยการนำข้อมูลจากเว็บไซต์ทั้งสิ้น 10 เว็บไซต์ จากคำค้นหา 5 คำที่แตกต่างกัน และถูกกำหนดคลาสโดยผู้วิจัยแล้ว เพื่อทำการทดสอบความถูกต้องของการ

ใช้เว็บบริการ โดยทำการอัปโหลดชุดข้อมูลทดสอบไปยังเว็บบริการเพื่อทำการจำแนกเนื้อหาจากไฟล์ข้อมูลดังกล่าว

The screenshot displays the BigML web interface for configuring a new batch prediction. The top navigation bar includes the BigML logo, user status (0 tasks, FREE), and various menu options. The main interface is titled 'PROJECT: ATapC' and 'New Batch Prediction'. It shows two input fields: '50 Webpages' and 'test_set dataset'. Below these are summary statistics for each dataset, including size, fields, instances, and models. A 'Configure' section is visible, followed by a 'Preview of the prediction file' showing a sample of the output data in a table format. At the bottom, there is a 'Prediction name' field and 'Reset' and 'Predict' buttons.

ภาพที่ 27 แสดงหน้าเว็บบริการในส่วนของการทดสอบตัวจำแนกด้วยชุดทดสอบ

ซึ่งจากภาพที่ 27 ที่ได้แสดงการทดสอบการจำแนก โดยกำหนดชุดข้อมูลทดสอบจากไฟล์ข้อมูลทดสอบที่ได้ทำการอัปโหลดไป เพื่อให้ได้ผลลัพธ์ในการคาดคะเนออกมา ซึ่งผลลัพธ์นี้จะถูกนำไปใช้ในส่วนของการทดสอบและการวิเคราะห์ผลต่อไป

บทที่ 5

การทดสอบและการวิเคราะห์ผล

5.1. วัตถุประสงค์ของการทดสอบ

จุดประสงค์ของการทดสอบ เพื่อสนับสนุนแนวทางในการสร้างเครื่องมือท่องเที่ยวเว็บไซต์บนอุปกรณ์เคลื่อนที่สำหรับคนพิการทางสายตาที่ได้ออกแบบ และพัฒนาเครื่องมือที่สนับสนุนแนวทางในบทที่ 4 โดยเนื้อหาจะประกอบไปด้วยการทดสอบระบบ ตั้งแต่การค้นหารายการเว็บไซต์ การสร้างไฟล์ข้อมูลเรียนรู้ ผลลัพธ์ของการค้นหาเนื้อหาหลัก การเปรียบเทียบการใช้คำค้นหาเป็นส่วนหนึ่งของคุณสมบัติในการค้นหาเนื้อหา ตลอดจนผลลัพธ์จากการใช้งานเครื่องมือที่จะสามารถลดจำนวนขั้นตอนในการเข้าถึงเนื้อหาของเว็บไซต์ได้

5.2. การทดสอบระบบ

การทดสอบเครื่องมือจะใช้วิธีทดสอบแบบปิด (Black-Box Testing) เพื่อต้องการที่จะทดสอบการทำงานของเครื่องมือให้เป็นไปตามความต้องการของฟังก์ชัน และมีการรองรับการทำงานของเครื่องมืออ่านหน้าจอ ซึ่งมีการทดสอบแบ่งออกเป็นดังต่อไปนี้

1) ทดสอบการค้นหารายการเว็บไซต์ที่เกี่ยวข้องด้วยคำค้นหาจากผู้ใช้

การทดสอบการค้นหารายการเว็บไซต์ มีจุดประสงค์เพื่อทดสอบผลลัพธ์การค้นหาด้วยคำค้นหาที่หลากหลายแบบที่อาจจะได้จากผู้ใช้ เช่น การค้นหาด้วยคำหลายคำ หรือการค้นหาด้วยคำว่างเปล่าในกล่องป้อนคำค้นหา ส่วนของเครื่องมือจะต้องสามารถตอบรับได้อย่างถูกต้อง เป็นต้น

ตารางที่ 3 ทดสอบการค้นหารายการเว็บไซต์ที่เกี่ยวข้องด้วยคำค้นหาจากผู้ใช้

ลำดับ	คำอธิบาย	ผลที่คาดหวัง	ผลลัพธ์
TC01	ทำการค้นหาโดยไม่ระบุคำค้นหาใดๆ	ระบบจะต้องแสดงการแจ้งเตือนเพื่อบ่งบอกให้ผู้ใช้ทราบว่าไม่สามารถดำเนินการได้	ถูกต้อง
TC02	ทำการค้นหาโดยระบุคำค้นหาจำนวน 1 คำ	ระบบจะต้องสร้างที่อยู่เว็บไซต์ได้อย่างถูกต้องและสามารถนำไปทดลองในการเปิดด้วยเว็บเบราว์เซอร์ได้	ถูกต้อง

TC03	ทำการค้นหาโดยระบุค่าค้นหาจำนวนมากกว่า 1 ค่า	ระบบจะต้องสร้างที่อยู่เว็บไซต์ได้อย่างถูกต้องและสามารถนำไปทดลองในการเปิดด้วยเว็บเบราว์เซอร์ได้	ถูกต้อง
TC04	ทำการค้นหาโดยระบุค่าค้นหา	ระบบจะต้องแสดงผลลัพธ์ได้ตรงกับการใช้เว็บไซต์ในการค้นหา	ถูกต้อง

2) ทดสอบการแปลงเอกสารเว็บไซต์ให้แสดงผลบนหน้าจอได้อย่างถูกต้อง

เป็นการทดสอบเพื่อตรวจสอบความถูกต้องของการแปลงเอกสารเว็บไซต์ที่ได้จากการส่งคำร้องขอไปยังเว็บไซต์กูเกิลเพื่อให้ได้รายการเว็บไซต์ที่เกี่ยวข้อง และทำการแปลงเป็นรายการเว็บไซต์ที่อยู่ในลักษณะที่ถูกต้องบนหน้าจอของเครื่องมือ

ตารางที่ 4 ทดสอบการแปลงเอกสารเว็บไซต์ให้แสดงผลบนหน้าจอได้อย่างถูกต้อง

ลำดับ	คำอธิบาย	ผลที่คาดหวัง	ผลลัพธ์
TC05	ผู้ใช้ป้อนค่าค้นหาและกดค้นหา	ระบบจะต้องสร้างรายการเว็บไซต์จากค่าค้นหาและแสดงเป็นลักษณะลิสต์ของหัวข้อได้อย่างถูกต้อง	ถูกต้อง
TC06	ผู้ใช้กดปุ่มย้อนกลับ	ระบบจะต้องเปลี่ยนหน้าจอกลับไปยังหน้าค้นหาได้	ถูกต้อง
TC07	ผู้ใช้ทำการเลือกเว็บไซต์ที่ต้องการให้ค้นหาเนื้อหาหลัก	ระบบจะต้องเปลี่ยนหน้าจอไปยังหน้าแสดงเนื้อหาหลักได้	ถูกต้อง

3) ทดสอบการสร้างไฟล์ข้อมูลจากที่อยู่เว็บไซต์

เป็นการทดสอบเพื่อตรวจสอบแปลงเอกสารของหน้าเว็บไซต์ที่ได้จากการเลือกของผู้ใช้ และทำการคำนวณเพื่อหาคุณลักษณะแบบต้นเพื่อใช้ในการสร้างตัวจำแนกในการคัดแยกเนื้อหาหลักบนหน้าเว็บไซต์ได้อย่างถูกต้อง อีกทั้งเพื่อเป็นการทดสอบการอัปโหลดไฟล์ข้อมูลที่จะถูกส่งไปยังเว็บบริการ

ตารางที่ 5 ทดสอบการสร้างไฟล์ข้อมูลจากที่อยู่เว็บไซต์

ลำดับ	คำอธิบาย	ผลที่คาดหวัง	ผลลัพธ์
TC08	สร้างไฟล์ข้อมูลคุณสมบัติแบบต้นจากเว็บไซต์ลำดับที่ 1	ระบบจะต้องสร้างไฟล์ข้อมูลที่อยู่ในรูปแบบที่มีการใช้จุลภาค	ถูกต้อง
TC09	สร้างไฟล์ข้อมูลเพื่อตรวจเนื้อหาจากเว็บไซต์ลำดับที่ 1	ระบบจะต้องสร้างไฟล์ข้อมูลที่มีเนื้อหาบรรจุอยู่และมีเลขเอกลักษณ์ (ID) สำหรับใช้สร้างไฟล์ข้อมูลเรียนรู้	ถูกต้อง
TC10	สร้างไฟล์ข้อมูลคุณสมบัติแบบต้นจากเว็บไซต์ลำดับที่ 2	ระบบจะต้องสร้างไฟล์ข้อมูลที่อยู่ในรูปแบบที่มีการใช้จุลภาค	ถูกต้อง
TC11	สร้างไฟล์ข้อมูลเพื่อตรวจเนื้อหาจากเว็บไซต์ลำดับที่ 2	ระบบจะต้องสร้างไฟล์ข้อมูลที่มีเนื้อหาบรรจุอยู่และมีเลขเอกลักษณ์ (ID) สำหรับใช้สร้างไฟล์ข้อมูลเรียนรู้	ถูกต้อง
TC12	สร้างไฟล์ข้อมูลคุณสมบัติแบบต้นจากเว็บไซต์ลำดับที่ 3	ระบบจะต้องสร้างไฟล์ข้อมูลที่อยู่ในรูปแบบที่มีการใช้จุลภาค	ถูกต้อง

ลำดับ	คำอธิบาย	ผลที่คาดหวัง	ผลลัพธ์
TC13	สร้างไฟล์ข้อมูลเพื่อตรวจเนื้อหาจากเว็บไซต์ลำดับที่ 3	ระบบจะต้องสร้างไฟล์ข้อมูลที่มีเนื้อหาบรรจุอยู่และมีเลขเอกลักษณ์ (ID) สำหรับใช้สร้างไฟล์ข้อมูลเรียนรู้	ถูกต้อง
TC14	สร้างไฟล์ข้อมูลคุณสมบัติแบบต้นจากเว็บไซต์ลำดับที่ 4	ระบบจะต้องสร้างไฟล์ข้อมูลที่อยู่ในรูปแบบที่มีการใช้จุลภาค	ถูกต้อง
TC15	สร้างไฟล์ข้อมูลเพื่อตรวจเนื้อหาจากเว็บไซต์ลำดับที่ 4	ระบบจะต้องสร้างไฟล์ข้อมูลที่มีเนื้อหาบรรจุอยู่และมีเลขเอกลักษณ์ (ID) สำหรับใช้สร้างไฟล์ข้อมูลเรียนรู้	ถูกต้อง
TC16	สร้างไฟล์ข้อมูลคุณสมบัติแบบต้นจากเว็บไซต์ลำดับที่ 5	ระบบจะต้องสร้างไฟล์ข้อมูลที่อยู่ในรูปแบบที่มีการใช้จุลภาค	ถูกต้อง
TC17	สร้างไฟล์ข้อมูลเพื่อตรวจเนื้อหาจากเว็บไซต์ลำดับที่ 5	ระบบจะต้องสร้างไฟล์ข้อมูลที่มีเนื้อหาบรรจุอยู่และมีเลขเอกลักษณ์ (ID) สำหรับใช้สร้างไฟล์ข้อมูลเรียนรู้	ถูกต้อง
TC18	สร้างไฟล์ข้อมูลจากเว็บไซต์ที่ค้นหาด้วยคำค้นหาภาษาอังกฤษ	ระบบจะต้องสร้างไฟล์ข้อมูลที่มีเนื้อหาได้อย่างถูกต้อง	ถูกต้อง
TC19	สร้างไฟล์ข้อมูลจากเว็บไซต์ที่ค้นหาด้วยคำค้นหาภาษาไทย	ระบบจะต้องสร้างไฟล์ข้อมูลที่มีเนื้อหาได้อย่างถูกต้อง	ไม่ถูก

ลำดับ	คำอธิบาย	ผลที่คาดหวัง	ผลลัพธ์
TC20	สร้างไฟล์ข้อมูลจากเว็บไซต์ที่มีการเข้ารหัสตัวอักษรแบบ UTF-8	ระบบจะต้องสร้างไฟล์ข้อมูลที่มีเนื้อหาถูกต้องไม่ผิดเพี้ยน	ถูกต้อง
TC21	สร้างไฟล์ข้อมูลจากเว็บไซต์ที่มีการเข้ารหัสตัวอักษรแบบ TIS-620	ระบบจะต้องสร้างไฟล์ข้อมูลที่มีเนื้อหาถูกต้องไม่ผิดเพี้ยน	ถูกต้อง
TC22	สร้างไฟล์ข้อมูลจากเว็บไซต์ที่มีการเข้ารหัสตัวอักษรแบบอื่นๆ	ระบบจะต้องสร้างไฟล์ข้อมูลที่มีเนื้อหาถูกต้องไม่ผิดเพี้ยน	ไม่ถูก

4) ทดสอบการเชื่อมต่อกับเว็บบริการในการค้นหาเนื้อหาหลักของเว็บไซต์

เป็นการทดสอบเพื่อตรวจสอบการแสดงผลหลังจากที่ได้ทำการวิเคราะห์เพื่อจำแนกหาเนื้อหาหลักบนหน้าเว็บไซต์จากเว็บบริการแล้ว ซึ่งเว็บบริการจะทำการส่งผลลัพธ์การวิเคราะห์กลับมา ทำให้การทดสอบจะต้องตรวจสอบว่าผลลัพธ์ที่ได้ตรงตามที่ต้องการหรือไม่

ตารางที่ 6 ทดสอบการเชื่อมต่อกับเว็บบริการในการค้นหาเนื้อหาหลักของเว็บไซต์

ลำดับ	คำอธิบาย	ผลที่คาดหวัง	ผลลัพธ์
TC23	สร้างตัวจำแนกบนเว็บบริการ	ระบบจะต้องสามารถอัปโหลดไฟล์ข้อมูลคุณลักษณะแบบต้นไปยังเว็บบริการได้	ถูกต้อง
TC24	จำแนกเนื้อหาหลักของหน้าเว็บไซต์	ระบบจะต้องสามารถส่งคำร้องขอไปยังเว็บบริการเพื่อประมวลผลในการใช้จำแนกเนื้อหาได้	ถูกต้อง

ลำดับ	คำอธิบาย	ผลที่คาดหวัง	ผลลัพธ์
TC25	ส่วนของการแสดงเนื้อหาหลัก	ระบบจะต้องสามารถรับผลลัพธ์การ จำแนกจากเว็บบริการและนำ ส่วนประกอบที่เป็นเนื้อหาหลักมา แสดงได้อย่างถูกต้อง	ถูกต้อง

5) ทดสอบการสนับสนุนเครื่องมืออ่านหน้าจอ

เป็นการทดสอบเพื่อให้แน่ใจว่าทุกส่วนประกอบบนหน้าจอของเครื่องมือ และลำดับ
การใช้งานของเครื่องมือมีการรองรับสนับสนุนเครื่องมืออ่านหน้าจอที่เหมาะสมในการถูกใช้
โดยผู้ใช้ที่ไม่สามารถมองเห็นหน้าจอได้

ตารางที่ 7 ทดสอบการสนับสนุนเครื่องมืออ่านหน้าจอ

ลำดับ	คำอธิบาย	ผลที่คาดหวัง	ผลลัพธ์
TC26	ส่วนของหน้าการค้นหา	เครื่องมืออ่านหน้าจอสามารถอ่าน ส่วนประกอบบนหน้าจอและ ช่วยเหลือในการป้อนคำค้นหาได้	ถูกต้อง
TC27	ส่วนของหน้าผลลัพธ์การค้นหา	เครื่องมืออ่านหน้าจอสามารถอ่าน ส่วนประกอบบนหน้าจอและ ช่วยเหลือให้ผู้พิการสามารถเข้าใจ หัวข้อแต่ละเว็บไซต์ได้	ถูกต้อง
TC28	ส่วนของหน้าผลลัพธ์การค้นหา	เครื่องมืออ่านหน้าจอสามารถอ่าน ส่วนประกอบบนหน้าจอและ ช่วยเหลือให้ผู้พิการสามารถเข้าใจ รายละเอียดโดยย่อแต่ละเว็บไซต์ได้	ถูกต้อง
TC29	ส่วนของหน้าผลลัพธ์การค้นหา	เครื่องมืออ่านหน้าจอสามารถ ช่วยเหลือให้ผู้ใช้สามารถย้อนกลับ ไปยังหน้าค้นหาได้	ถูกต้อง

ลำดับ	คำอธิบาย	ผลที่คาดหวัง	ผลลัพธ์
TC30	ส่วนของหน้าแสดงเนื้อหาหลัก	เครื่องมืออ่านหน้าจอจะต้องสามารถบ่งบอกสถานะของการทำงานโปรแกรมได้	ถูกต้อง
TC31	ส่วนของหน้าแสดงเนื้อหาหลัก	เครื่องมืออ่านหน้าจอสามารถอ่านเนื้อหาหลักบนหน้าจอ	ถูกต้อง
TC32	ส่วนของหน้าแสดงเนื้อหาหลัก	เครื่องมืออ่านหน้าจอสามารถช่วยเหลือให้ผู้ใช้กลับไปยังหน้าผลลัพธ์การค้นหาได้	ถูกต้อง

6) ทดสอบการนำคำค้นหาใช้ในการค้นหาเนื้อหา

เป็นการทดสอบเพื่อแสดงให้เห็นถึงประโยชน์ของการนำคำค้นหาจากผู้ใช้มาช่วยในการค้นหาเนื้อหาหลักบนหน้าเว็บไซต์ โดยใช้โปรแกรม Weka เข้ามาช่วยในการแสดงประสิทธิภาพของการค้นหา โดยการทดสอบนี้จะนำไฟล์ข้อมูลที่ได้เตรียมไว้มาทำการสร้างตัวจำแนกเนื้อหาด้วยวิธีป่าแบบสุ่ม มีการกำหนดค่าให้มีจำนวนต้นไม้ทั้งหมด 100 ต้น มีความลึกสูงสุด (Max Depth) อยู่ที่ 6 ระดับ จำนวนรอบ (Iteration) ที่ 500 รอบ และทำการทดสอบแบบ 10-fold Cross Validation ซึ่งมีผลลัพธ์ดังรูปภาพต่อไปนี้

Classifier: RandomForest -P 100 -print -I 500 -num-slots 1 -K 1 -M 1.0 -V 0.001 -S 1 -depth 6

Test options:
 Use training set
 Supplied test set (Set...)
 Cross-validation Folds: 10
 Percentage split %: 66
 More options...

(Nom) class: [Dropdown]
 Start Stop

Result list (right-click for options):
 20:40:11 - trees.RandomForest
 20:42:23 - trees.RandomForest
 20:44:04 - trees.RandomForest
 20:50:09 - trees.RandomForest
 20:51:00 - trees.RandomForest
 20:52:14 - trees.RandomForest

Classifier output:
 === Stratified cross-validation ===
 === Summary ===
 Correctly Classified Instances 20753.407 88.6254 %
 Incorrectly Classified Instances 2663.593 11.3746 %
 Kappa statistic 0.7725
 Mean absolute error 0.1794
 Root mean squared error 0.2793
 Relative absolute error 35.8793 %
 Root relative squared error 55.8522 %
 Total Number of Instances 23417

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC
	0.933	0.161	0.853	0.933	0.891	0.776	0.9
	0.839	0.067	0.926	0.839	0.881	0.776	0.9
Weighted Avg.	0.886	0.114	0.890	0.886	0.886	0.776	0.9

=== Confusion Matrix ===

a	b	<-- classified as
10927	782	a = FALSE
1882	9827	b = TRUE

Status: OK Log x 0

ภาพที่ 28 ผลลัพธ์จากโปรแกรม Weka ด้วยไฟล์ข้อมูลที่ไม่มีการใช้คำค้นหา

Classifier: RandomForest -P 100 -print -I 500 -num-slots 1 -K 1 -M 1.0 -V 0.001 -S 1 -depth 5

Test options:
 Use training set
 Supplied test set (Set...)
 Cross-validation Folds: 10
 Percentage split %: 66
 More options...

(Nom) class: [Dropdown]
 Start Stop

Result list (right-click for options):
 20:06:30 - trees.RandomForest
 20:30:29 - trees.RandomForest
 20:31:34 - trees.RandomForest
 20:32:40 - trees.RandomForest
 20:33:39 - trees.RandomForest
 20:34:32 - trees.RandomForest
 20:35:25 - trees.RandomForest

Classifier output:
 === Stratified cross-validation ===
 === Summary ===
 Correctly Classified Instances 21447.3075 91.5886 %
 Incorrectly Classified Instances 1969.6925 8.4114 %
 Kappa statistic 0.8318
 Mean absolute error 0.1956
 Root mean squared error 0.2737
 Relative absolute error 39.1158 %
 Root relative squared error 54.7269 %
 Total Number of Instances 23417

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC
	0.939	0.107	0.898	0.939	0.918	0.833	0.9
	0.893	0.061	0.936	0.893	0.914	0.833	0.9
Weighted Avg.	0.916	0.084	0.917	0.916	0.916	0.833	0.9

=== Confusion Matrix ===

a	b	<-- classified as
10993.29	715.21	a = FALSE
1254.48	10454.02	b = TRUE

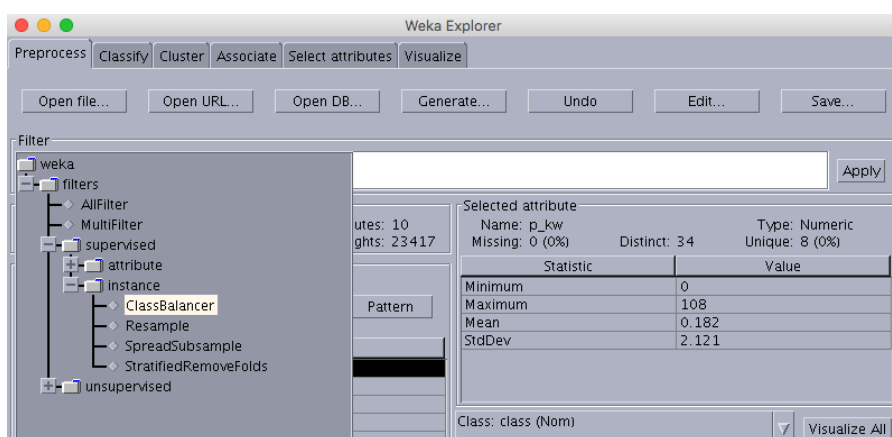
Status: OK Log x 0

ภาพที่ 29 ผลลัพธ์จากโปรแกรม Weka ด้วยไฟล์ข้อมูลที่มีการใช้คำค้นหาเป็นส่วนประกอบ

ซึ่งจะเห็นได้ว่าผลลัพธ์ของการสร้างวิธีการป่าแบบสุ่มที่ประกอบไปด้วยการเพิ่มคำค้นหานั้น จะช่วยให้สามารถแยกส่วนที่เป็นเนื้อหาและไม่ใช่นี้อาได้แม่นยำขึ้นจากเดิมถึงร้อยละ 2.96

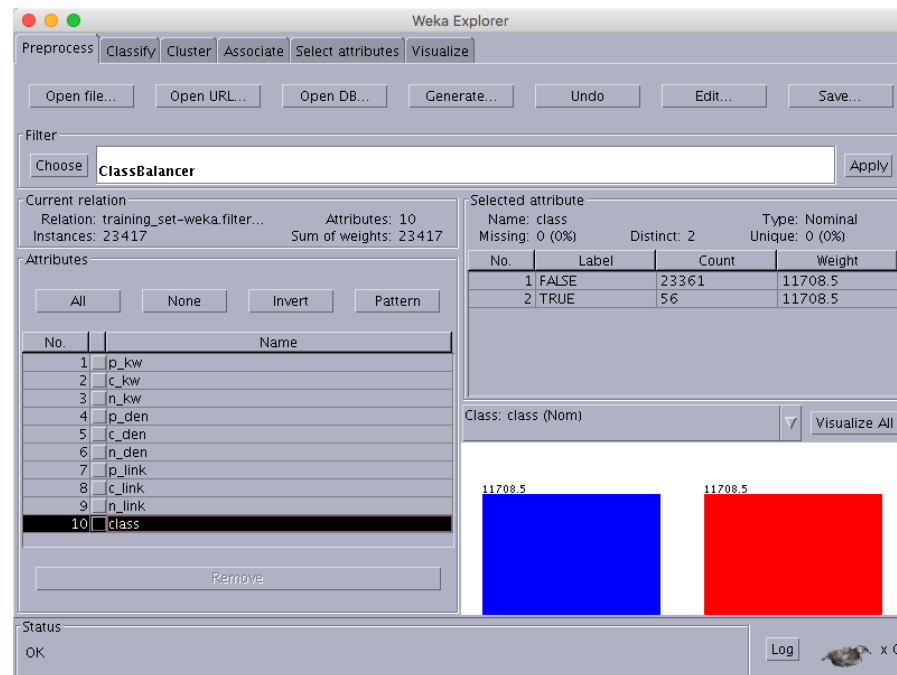
7) ทดสอบเพื่อเปรียบเทียบความถูกต้องในการค้นหาเนื้อหาหลักด้วยวิธีการป่าแบบสุ่ม กับวิธีการแบบอื่นๆ

เป็นการทดสอบเพื่อเปรียบเทียบการค้นหาเนื้อหาหลักด้วย 3 วิธีการที่ต่างกันอย่างสิ้นเชิง โดยผู้วิจัย จะทำการนำชุดข้อมูลเรียนรู้ (Training Set) เดียวกันมาใช้ ร่วมกับการทดสอบแบบ 10-fold cross validation ซึ่งผู้วิจัยได้ทำการคัดเลือกมา 2 วิธีการมาเปรียบเทียบได้แก่ วิธีแบบต้นไม้ตัดสินใจ (Decision Tree) และวิธีแบบโครงข่ายประสาท (Neural Network) โดยมีวิธีการทดสอบดังนี้



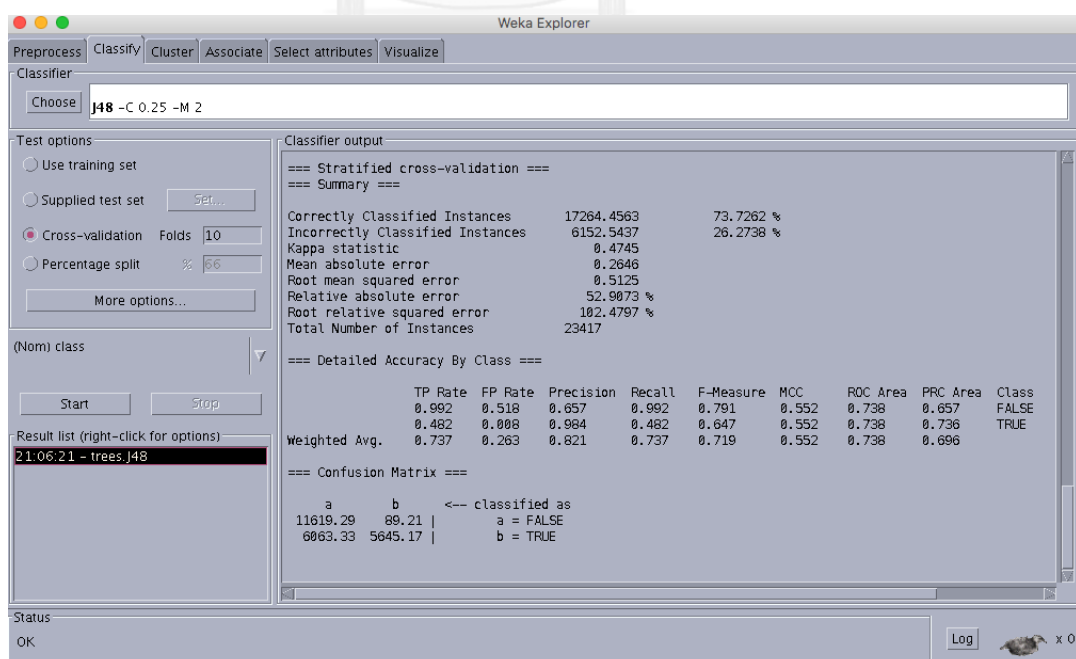
ภาพที่ 30 การเพิ่มตัวกรองในโปรแกรม Weka

จากภาพที่ 30 ในขั้นต้นผู้วิจัยจะทำการนำเข้าไฟล์ข้อมูลเรียนรู้ดังที่ได้จากหัวข้อ 3.3 เข้าสู่เครื่องมือ Weka Explorer โดยก่อนที่จะนำข้อมูลเหล่านี้ไปสร้างตัวจำแนกจะต้องมีการใช้ตัวกรองข้อมูล ClassBalancer เพื่อสร้างสมดุลน้ำหนักของคลาสที่เป็นเนื้อหาและไม่เป็นเนื้อหา เนื่องจากจำนวนข้อมูลที่เป็นเนื้อหาหลักมีจำนวนน้อยกว่ามากเมื่อเปรียบเทียบกับส่วนอื่นๆ ซึ่งอาจจะเกิดข้อมูลที่มือคุดเกิดขึ้นได้ เมื่อทำการเลือกตัวกรองแบบ ClassBalancer แล้วให้กดปุ่ม Apply เพื่อใช้งานตัวกรองข้อมูลดังกล่าวกับชุดข้อมูลที่นำเข้า



ภาพที่ 31 ผลลัพธ์การใช้ตัวกรองข้อมูล ClassBalancer

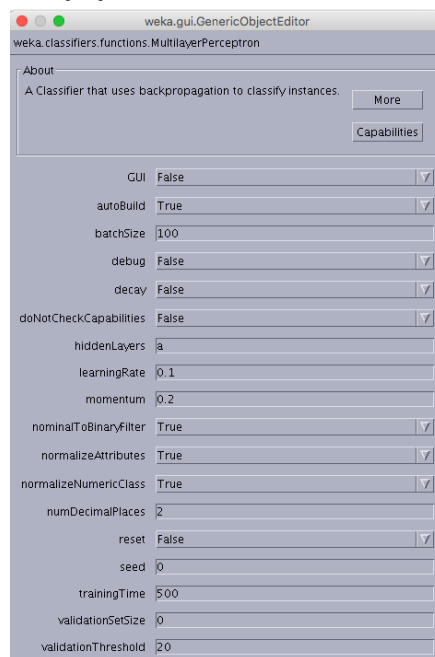
หลังจากการใช้ตัวกรองเพื่อสร้างสมดุลแล้ว จะได้ผลลัพธ์ดังรูปภาพที่ 31 ซึ่งแสดงให้เห็นถึงน้ำหนักของคลาส (Weight) ที่เท่ากันระหว่างคลาส TRUE และ FALSE จากนั้นผู้วิจัยจะทำการสร้างตัวจำแนกเนื้อหาจากชุดข้อมูลเรียนรู้ที่ได้มา โดยในขั้นแรกผู้วิจัยจะทำการสร้างตัวจำแนกด้วยวิธีแบบต้นไม้ตัดสินใจ โดยการเลือกตัวจำแนก (Classifier) เป็นประเภท J48 และกำหนดค่าแบบมาตรฐานเริ่มต้น



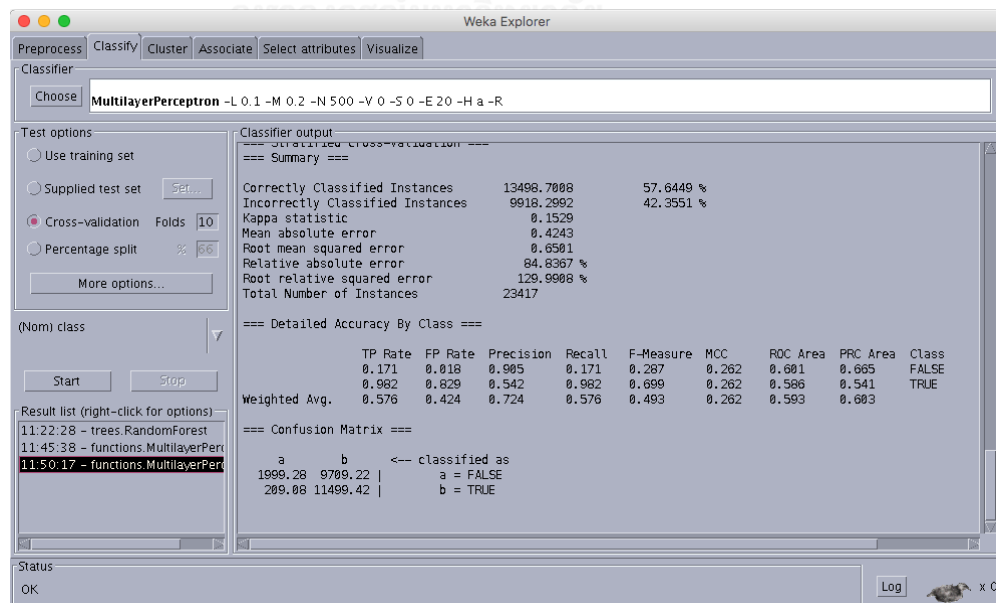
ภาพที่ 32 ผลลัพธ์การทดสอบด้วยตัวจำแนกแบบต้นไม้ตัดสินใจ

จากภาพที่ 32 ได้แสดงให้เห็นผลลัพธ์จากการทดสอบตัวจำแนกแบบต้นไม้ตัดสินใจด้วยวิธี 10-fold cross validation ทำให้ได้ผลลัพธ์ความถูกต้องในการจำแนกระหว่างส่วนประกอบอื่นๆ กับ ส่วนประกอบที่เป็นเนื้อหาหลักอยู่ที่ 73.72%

ในการสร้างตัวจำแนกแบบที่สอง คือ การสร้างด้วยวิธีแบบโครงข่ายประสาท (Neural Network) ด้วยการใช้ Multi-layer perceptron โดยมีการกำหนดค่าพารามิเตอร์ดังภาพที่ 33 โดยผู้วิจัยได้กำหนดให้มีระดับการเรียนรู้อยู่ที่ 0.1 และจำนวนการฝึกฝนอยู่ที่ 500 ครั้ง



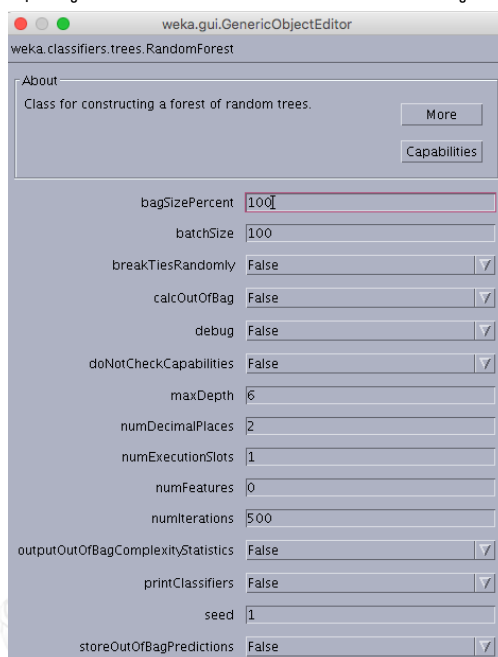
ภาพที่ 33 แสดงการกำหนดค่าพารามิเตอร์ในการสร้างตัวจำแนกแบบโครงข่ายประสาท



ภาพที่ 34 ผลลัพธ์การทดสอบด้วยตัวจำแนกแบบโครงข่ายประสาท

จากภาพที่ 34 ได้แสดงให้เห็นผลลัพธ์จากการทดสอบตัวจำแนกแบบโครงข่ายประสาทด้วยวิธี 10-fold cross validation เช่นเดียวกัน ทำให้ได้ผลลัพธ์ความถูกต้องในการจำแนกระหว่างส่วนประกอบอื่นๆ กับส่วนประกอบที่เป็นเนื้อหาหลักอยู่ที่ 57.84%

และในการสร้างตัวจำแนกแบบสุดท้าย คือ การสร้างด้วยวิธีแบบป่าแบบสุ่ม (Random Forest) โดยมีการกำหนดค่าพารามิเตอร์ดังภาพที่ 35 โดยผู้วิจัยได้กำหนดให้มีการสร้างต้นไม้ตัดสินใจจำนวน 100 ต้นกำหนดความลึกสูงสุดอยู่ที่ 6 ระดับและจำนวนการฝึกฝนอยู่ที่ 500 ครั้ง



ภาพที่ 35 แสดงการกำหนดค่าพารามิเตอร์ในการสร้างตัวจำแนกด้วยวิธีป่าแบบสุ่ม

Classifier output

```

=== Classifier Cross-Validation ===
=== Summary ===
Correctly Classified Instances 20541.7401      87.7215 %
Incorrectly Classified Instances 2875.2599      12.2785 %
Kappa statistic 0.7544
Mean absolute error 0.1489
Root mean squared error 0.2906
Relative absolute error 29.767 %
Root relative squared error 58.1106 %
Total Number of Instances 23417

=== Detailed Accuracy By Class ===
          TP Rate  FP Rate  Precision  Recall  F-Measure  MCC  ROC Area  PRC Area  Class
Weighted Avg.  0.877  0.123  0.886  0.877  0.877  0.763  0.972  0.968  TRUE
0.951  0.196  0.829  0.951  0.886  0.763  0.972  0.976  FALSE
0.884  0.049  0.942  0.884  0.867  0.763  0.972  0.968  TRUE
  
```

=== Confusion Matrix ===

	a	b	←- classified as
a	11133.12	575.38	a = FALSE
b	2299.88	9408.62	b = TRUE

ภาพที่ 36 ผลลัพธ์การทดสอบตัวจำแนกที่สร้างด้วยวิธีป่าแบบสุ่ม

จากภาพที่ 36 ได้แสดงให้เห็นผลลัพธ์จากการทดสอบตัวจำแนกที่สร้างด้วยวิธีป่าแบบสุ่มด้วยการใช้วิธี 10-fold cross validation เช่นเดียวกัน ซึ่งผลลัพธ์ความถูกต้องในการจำแนกระหว่างส่วนประกอบอื่นๆ กับส่วนประกอบที่เป็นเนื้อหาหลักอยู่ที่ 87.72% ซึ่งมีความถูกต้องมากที่สุดในการเปรียบเทียบกับวิธีการอื่นๆ ที่กล่าวมาข้างต้น

โดยจากการทดสอบข้างต้นทำให้สามารถสรุปผลได้ว่า เมื่อใช้ชุดข้อมูลเดียวกันโดยมีการใช้ตัวกรอง ClassBalancer เพื่อทำการสร้างสมดุลให้กับน้ำหนักของคลาสในชุดข้อมูล แล้วจึงนำไปใช้ทดสอบกับตัวจำแนกทั้งสามวิธี ผลลัพธ์ที่ได้คือ ตัวจำแนกที่ดีที่สุดคือวิธีป่าแบบสุ่ม (Random Forest) ซึ่งมีความถูกต้องในการจำแนกสูงที่สุดอยู่ที่ 87.72% ลำดับถัดมาคือตัวจำแนกแบบต้นไม้ตัดสินใจ (Decision Tree : J48) ซึ่งมีความถูกต้องอยู่ที่ 73.72% และลำดับสุดท้ายคือการใช้โครงข่ายประสาทด้วยวิธีแบบเพอร์เซ็ปตรอนหลายชั้น (Multi-layer perceptron) โดยได้ผลลัพธ์ความถูกต้องในการจำแนกอยู่ที่ 57.84%

8) ทดสอบขั้นตอนที่ใช้ในการเข้าถึงเนื้อหาเว็บไซต์โดยใช้วิธีการป่าแบบสุ่ม

เป็นการทดสอบเพื่อเปรียบเทียบการเข้าถึงเนื้อหาเว็บไซต์ด้วยวิธีการใช้เครื่องมือท่องเว็บ Safari กับการใช้เครื่องมือของผู้วิจัยที่ใช้วิธีป่าแบบสุ่มเพื่อตรวจหาเนื้อหาหลักสำคัญบนหน้าเว็บไซต์ โดยเปรียบเทียบจากจำนวนขั้นตอนหลังจากเปิดหน้าเว็บไซต์ที่ต้องการของเครื่องมือท่องเว็บ Safari กับหน้าแสดงเนื้อหาหลักในเครื่องมือของผู้วิจัยด้วยที่อยู่เว็บไซต์เดียวกัน และนับจำนวนครั้งในการกระทำเพื่อเลื่อนตัวเลือกของเครื่องมืออ่านหน้าจอ ดังที่แสดงในรูปภาพที่ 37



ภาพที่ 37 หน้าจอเริ่มต้นการทดสอบโดยเครื่องมือท่องเว็บ Safari (ซ้าย)

และ เครื่องมือของผู้วิจัย (ขวา)

โดยในการทดสอบนี้จะทำการสุ่มหาเว็บไซต์ตัวอย่าง ที่จะนำมาทดสอบจำนวนทั้งสิ้น 20 เว็บไซต์ที่มีความแตกต่างกัน และได้ผลการทดลองดังนี้

ตารางที่ 8 แสดงการเปรียบเทียบขั้นตอนในการเข้าถึงส่วนที่เป็นเนื้อหาหลักของหน้าเว็บไซต์

ลำดับ	รายชื่อเว็บไซต์	จำนวนขั้นตอนโดยใช้ เครื่องมือท่องเว็บ Safari	จำนวนขั้นตอนโดยใช้ เครื่องมือในงานวิจัย
1	ครอบครัวข่าว 3	29	6
2	ช่อง 7 สี	32	7
3	VOA	12	5
4	Post Today	16	4
5	สยามกีฬา	24	3
6	คม ชัด ลึก	10	3
7	Thai PR	17	7
8	Sanook	15	5
9	INN News	28	4
10	สยามโฟน	39	3
11	Bear Tai	43	5
12	Now 26	33	6
13	ARip	45	6
14	ประชาชาติ	17	2
15	ไทยรัฐ	41	7
16	ฐานเศรษฐกิจ	?	11
17	ผู้จัดการ	58	?
18	TechMoBlog	17	6
19	MThai	26	5
20	Spring News	29	7

หมายเหตุ: เครื่องหมายคำถามหรือปรัศนี มีความหมายว่าไม่สามารถหาเนื้อหาหลักได้บนหน้าเว็บไซต์

จากผลการทดสอบในตารางที่ 8 ได้แสดงให้เห็นถึงจำนวนขั้นตอนที่ใช้ในการเข้าถึงเนื้อหาด้วยเครื่องมือท่องเว็บ Safari กับเครื่องมือท่องเว็บของผู้วิจัย ซึ่งแสดงให้เห็นว่าเครื่องมือท่องเว็บของผู้วิจัยสามารถลดขั้นตอนการเข้าถึงเนื้อหาหลักจากวิธีเดิมซึ่งอยู่ที่ 27 ขั้นตอนโดยเฉลี่ย ลงมาเหลืออยู่ที่ 6 ขั้นตอนโดยเฉลี่ย ถึงแม้ว่าในบางเว็บไซต์ตัวอย่างเช่น VOA ได้มีการพัฒนาเพื่อรองรับ

การเข้าถึงโดยผู้พิการ ด้วยการสร้างลิงก์กระโดดไปส่วนเนื้อหาหลัก (Skip link) แต่ยังคงปัญหาว่าหลังจากที่ได้ใช้ลิงก์ดังกล่าวแล้ว ก็ยังพบอุปสรรคอื่นๆ เช่น ปุ่มที่ใช้ในการแบ่งปันข่าวนี้บนโซเชียลเน็ตเวิร์คต่างๆ ที่มีจำนวนอยู่หลายปุ่มก่อนที่จะถึงส่วนที่เป็นเนื้อหาหลักได้ เป็นต้น

จากผลการทดสอบดังกล่าวผู้วิจัยจึงสามารถสรุปผลได้ว่าเครื่องมือของผู้วิจัยทำให้การเข้าถึงเนื้อหาหลักเป็นไปได้อย่างรวดเร็วมากขึ้น โดยใช้การค้นหาเนื้อหาจากคุณสมบัติแบบต้นร่วมกับจำนวนคำค้นหา และวิธีการแบบปาสุ่ม

5.3. สรุปผลการทดสอบ

ผู้วิจัยสามารถสรุปผลการทดสอบตามวัตถุประสงค์ของการทดสอบได้ดังต่อไปนี้

- 1) การค้นหารายการเว็บไซต์ที่เกี่ยวข้องด้วยคำค้นหาจากผู้ใช้งานสามารถสร้างผลลัพธ์การค้นหาได้อย่างถูกต้อง
- 2) ผลลัพธ์จากการค้นหาสามารถนำมาทำการแปลงเอกสารเว็บไซต์ให้แสดงผลบนหน้าจอได้อย่างถูกต้อง
- 3) ตัวเครื่องมือสามารถสร้างข้อมูลจากที่อยู่เว็บไซต์ที่ถูกเลือกจากผู้ใช้ได้อย่างถูกต้อง และมีขนาดที่ไม่ใหญ่มากเกินไป จึงทำให้สามารถอัปโหลดไปที่เว็บบริการได้ด้วยความเร็ว
- 4) ในการค้นหาเนื้อหาหลักของหน้าเว็บไซต์ จากการทดสอบทำให้พบว่าการนำคำค้นหาจากผู้ใช้งานมาช่วยในการค้นหาจะทำให้การค้นหามีความแม่นยำมากยิ่งขึ้น แต่ยังคงพบปัญหาหลักในการแสดงผลในบางเว็บไซต์เนื่องจากในบางเว็บไซต์ได้มีการกำหนดการเข้ารหัสของตัวอักษรที่ทำให้การส่งคำร้องขอจากตัวเครื่องมือทำให้เกิดเนื้อหาที่ไม่สามารถเข้าใจได้ มีลักษณะผิดเพี้ยนของตัวอักษร ซึ่งเกิดจากการที่ระบบไม่สามารถแปลงเอกสารให้อยู่ในรหัสแบบ UTF-8 ได้ จึงทำให้เกิดความผิดพลาดในการแสดงเนื้อหาจากหน้าเว็บไซต์เหล่านั้น
- 5) ในการทดสอบการสนับสนุนเครื่องมืออ่านหน้าจอ มีการกำหนดให้แต่ละส่วนประกอบบนหน้าจอให้ช่วยเหลือให้ผู้ที่มีความพิการทางสายตาให้สามารถใช้งานได้อย่างง่ายขึ้น
- 6) ในการทดสอบเพื่อวัดผลของการนำวิธีการปาสุ่มมาใช้ในการค้นหาเนื้อหา สามารถบ่งชี้ได้ว่าทำให้สามารถเข้าถึงเนื้อหาหลักของหน้าเว็บไซต์ได้อย่างรวดเร็วมากยิ่งขึ้น เนื่องมาจากการทำให้ผู้ใช้งานสามารถเข้าถึงเนื้อหาได้ทันที ในกรณีที่การจำแนกเนื้อหาเป็นไปอย่างถูกต้อง หรือในกรณีอื่นจะต้องอาศัยขั้นตอนเพิ่มโดยเฉลี่ยที่ 6 ขั้นตอนเพื่อไปยังจุดที่เป็นเนื้อหาหลัก เมื่อเปรียบเทียบกับกรเข้าหน้าเว็บไซต์ปกติโดยการใช้แอปพลิเคชัน Safari ที่ต้องใช้ขั้นตอนการเข้าถึงนับตั้งแต่เปิดหน้าเว็บไซต์จะอยู่ที่ 27 ขั้นตอนโดยเฉลี่ย

บทที่ 6

สรุปผลการวิจัย

ในบทนี้จะกล่าวถึงการสรุปผลงานวิจัย ข้อจำกัดของงานวิจัย และงานวิจัยที่สามารถพัฒนาต่อในอนาคตจากวิทยานิพนธ์ โดยแต่ละส่วนที่กล่าวมานั้นมีรายละเอียดดังต่อไปนี้

6.1. สรุปผลการวิจัย

เครื่องมือทอ้งเว็บไซต์นี้ได้ถูกออกแบบมาให้สามารถช่วยเหลือการเข้าถึงเนื้อหาของหน้าเว็บไซต์สำหรับผู้พิการทางสายตา ซึ่งการทำงานของเครื่องมืออาศัยเพียงคำค้นหาจากผู้ใช้ทำให้สามารถค้นหาเนื้อหาข่าวในทุกๆภาษาหรือหัวข้อข่าวที่ผู้ใช้สนใจได้ โดยทั่วไปการทอ้งเว็บไซต์ของผู้พิการทางสายตาตามักจะเต็มไปด้วยอุปสรรคในการเข้าถึงไม่ว่าจะเป็นการพัฒนาเว็บไซต์ที่ไม่ได้คำนึงถึงการใช้งานโดยผู้พิการทางสายตา จึงไม่ได้พัฒนาเว็บไซต์ตามข้อแนะนำของ WCAG 2.0 หรือในเว็บไซต์ที่พัฒนาด้วยข้อแนะนำเหล่านั้นแล้ว ก็ยังปรากฏงานวิจัยที่พิสูจน์ได้ว่ายังไม่สามารถทำให้ผู้พิการทางสายตาเข้าถึงเนื้อหาเหล่านั้นได้อย่างถูกต้อง อันเนื่องมาจากเนื้อหาที่ได้นั้นไม่ใช่เนื้อหาหลักบนหน้าเว็บไซต์ที่ผู้ใช้ต้องการ ผู้วิจัยจึงต้องการที่จะสร้างเครื่องมือการทอ้งเว็บไซต์ที่จะสามารถช่วยลดระยะเวลาและขั้นตอนในการเข้าถึงเนื้อหาบนหน้าเว็บไซต์ที่ซับซ้อนและเต็มไปด้วยสิ่งแปลกปลอมที่ไม่ใช่เนื้อหาหลัก ซึ่งจะช่วยให้ผู้พิการทางสายตาสามารถเข้าถึงได้ด้วยขั้นตอนที่น้อยกว่า นับตั้งแต่การเข้าสู่เครื่องมือในหน้าค้นหาที่ประกอบไปด้วยกล่องข้อมูลเพื่อใช้ในการป้อนคำค้นหาผ่านทางแป้นพิมพ์ที่มีการรองรับเพื่อผู้พิการทางสายตา หรือการใช้เครื่องมือในการแปลงเสียงให้เป็นตัวอักษร จากนั้นเมื่อผู้ใช้งานทำการกดค้นหาด้วยการใช้เครื่องมืออ่านหน้าจอด้วยความช่วยเหลือที่ออกแบบให้สามารถใช้งานได้โดยง่ายตายแล้วนั้น เครื่องมือจะทำการแสดงผลการค้นหาออกมาในรูปแบบรายการของเว็บไซต์ที่เกี่ยวข้องและทำให้ผู้ใช้สามารถไปสู่รายละเอียดของเว็บไซต์ที่ละเว็บได้โดยผ่านทางการใช้เครื่องมืออ่านหน้าจอในการอ่านหัวข้อของหน้าเว็บไซต์และเนื้อหาโดยย่อ ซึ่งจะช่วยให้ผู้ใช้สามารถตัดสินใจในการเว็บไซด์ได้ที่ต้องการได้อย่างถูกต้อง จากนั้นเมื่อผู้ใช้งานเลือกเว็บไซด์ที่ต้องการแล้ว เครื่องมือจะแสดงเนื้อหาที่ได้จากการจำแนกออกมาโดยวิธีการป่าแบบสุ่มและทำการอ่านให้ผู้ใช้สามารถเข้าใจถึงเนื้อหาบนเว็บไซด์นั้นๆ ได้อย่างรวดเร็วขึ้น

จากการทำการทดสอบและวิเคราะห์ผลลัพธ์ที่ได้จากบทที่ 5 ได้แสดงให้เห็นถึงการใช้งานเครื่องมือที่ประกอบไปด้วยตัวเครื่องมือบนอุปกรณ์เคลื่อนที่ ประมวลผลร่วมกับการใช้งานเว็บบริการเพื่อทำการจำแนกเนื้อหาหลักออกจากหน้าเว็บไซต์ และมีการสนับสนุนของเครื่องมืออ่านหน้าจอในทุกๆ ส่วนประกอบเพื่ออำนวยความสะดวกให้กับผู้พิการทางสายตาในการเข้าถึงหน้าเว็บไซต์ต่างๆ ซึ่งผู้วิจัย

ได้หวังเป็นอย่างยิ่งว่าเครื่องมือขั้นนี้จะสามารถนำไปพัฒนาและใช้งานต่อไปในอนาคตเพื่อเป็นประโยชน์แก่ผู้พิการทางสายตาที่ต้องการจะเข้าถึงเนื้อหาสาระอันเป็นประโยชน์

6.2. ข้อจำกัดของงานวิจัย

- 1) เครื่องมือสามารถใช้งานได้เต็มประสิทธิภาพบนระบบ iOS เวอร์ชัน 9 ขึ้นไป
- 2) ในอุปกรณ์จะต้องมีการติดตั้งเครื่องมืออ่านหน้าจอ VoiceOver ด้วยเพื่อที่จะทำให้ผู้ใช้ที่ไม่สามารถมองเห็นใช้งานได้
- 3) การแปลงเอกสารเว็บไซต์นั้น หน้าเว็บไซต์จะต้องมีการรองรับการแปลงเป็นรหัส UTF-8 ได้ในภาษาไทย มิเช่นนั้นการแสดงผลจะผิดเพี้ยน
- 4) การค้นหาเนื้อหาจะต้องต่อกับส่วนประสานของเว็บบริการ ซึ่งในต้นแบบนี้เว็บบริการจะถูกจำกัดสิทธิ์การใช้งานให้ได้ทีละหนึ่งคนและหนึ่งการทำงานเท่านั้น ไม่สามารถรองรับการทำงานแบบขนานได้
- 5) ผู้ใช้จะต้องป้อนคำค้นหาด้วยตัวเอง ซึ่งจะต้องอาศัยเครื่องมือเสริมเช่น แป้นพิมพ์ที่มีการสนับสนุนจากเครื่องมืออ่านหน้าจอ หรือการใช้งาน Speech-to-text ในการแปลงเสียงพูดให้เป็นตัวอักษร

6.3. งานวิจัยในอนาคต

- 1) เพิ่มความรวดเร็วในการเข้าถึงเนื้อหา ด้วยการสร้างส่วนแนะนำเว็บไซต์ที่ความเกี่ยวข้องกับคำค้นหาก่อนหน้านี้ ส่วนนี้จะช่วยทำให้ผู้ใช้สามารถเข้าถึงเว็บไซต์ที่สนใจได้ตั้งแต่นำค้นหา โดยไม่จำเป็นต้องทำการป้อนคำค้นหาใหม่อีกครั้ง
- 2) เครื่องมือควรมีการสนับสนุนการแปลงเอกสารจากการเข้ารหัสตัวอักษรในรหัสแบบอื่นๆได้อย่างสมบูรณ์แบบมากขึ้น
- 3) ควรพัฒนาเว็บบริการให้มีการให้บริการในการจำแนกเนื้อหาโดยตรง เพื่อหลีกเลี่ยงปัญหาในการทำงานหลายขั้นตอน เพราะเนื่องจากเว็บบริการที่ใช้ในงานวิจัยนี้มีขีดจำกัดในการใช้งาน และส่วนต่อประสานไม่ได้มีการพัฒนาต่อทำให้ยากต่อการใช้งาน และรับส่งข้อมูลระหว่างเครื่องมือและเว็บบริการ

รายการอ้างอิง

- [1] V. L. Hanson, "Progress on Website Accessibility?," *ACM Trans. Web*, vol. 7, p. 30, 2013.
- [2] C. Kohlschütter, "Boilerplate Detection using Shallow Text Features," presented at the WSDM'10, 2010.
- [3] L. Breiman, "Random Forests," *Machine Learning*, pp. 5-32, January 2001 2001.
- [4] P. P. Y. Lai, "Application of Content Adaptation in Web Accessibility for the Blind," *W4A2011 - Communications*, March 28-29, 2011 2011.
- [5] E. Lundgren, "Extracting news text from web pages: an application for the visually impaired," *PETRA '15*, July 01 - 03 2015 2015.
- [6] BigML. (2011, 23 July). *BigML is Machine Learning for everyone*. Available: <https://bigml.com/>



ภาคผนวก

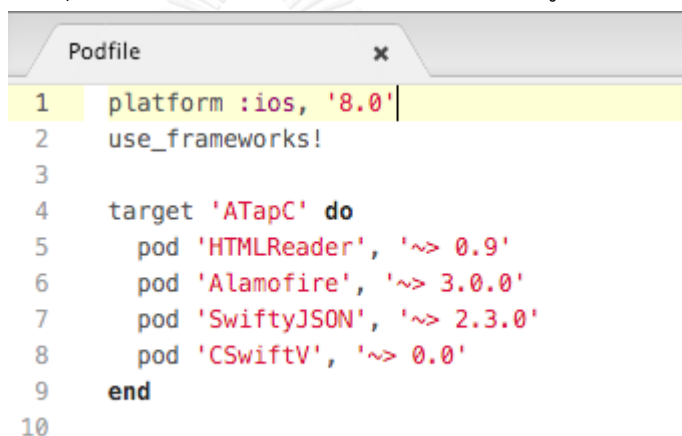
จุฬาลงกรณ์มหาวิทยาลัย
CHULALONGKORN UNIVERSITY

ภาคผนวก ก.

การติดตั้งส่วนเสริมโคโคพอด (Cocoa Pod), อลาโมไฟร์ (Alamofire) และตัวอ่านเอกสารเว็บไซต์ (HTML Reader)

ในงานวิจัยนี้ได้พัฒนาเครื่องมือด้วยภาษา Swift และใช้โปรแกรม Xcode ในการเขียนและคอมไพล์โปรแกรม ซึ่งจำเป็นที่จะต้องอาศัยส่วนเสริมจำนวนหนึ่งในการพัฒนาให้เครื่องมือสามารถใช้การติดต่อแบบ HTTP ได้อย่างสะดวกมากขึ้น โดยมีการติดตั้งดังต่อไปนี้

- 1) ติดตั้งตัวจัดการแพ็คเกจอัตโนมัติ Cocoa Pod ด้วยการใช้คำสั่งดังต่อไปนี้บนเทอร์มินอล
\$ sudo gem install cocoapods
- 2) ติดตั้งส่วนเสริมอลาโมไฟร์ (Alamofire) และตัวอ่านเอกสารเว็บไซต์ (HTML Reader) ด้วยการสร้างไฟล์ที่บรรจุส่วนเสริมที่ต้องการ โดยมีตัวอย่างไฟล์ดังในรูปภาพที่ 38



```

Podfile
1 platform :ios, '8.0'
2 use_frameworks!
3
4 target 'ATapC' do
5   pod 'HTMLReader', '~> 0.9'
6   pod 'Alamofire', '~> 3.0.0'
7   pod 'SwiftyJSON', '~> 2.3.0'
8   pod 'CSwiftV', '~> 0.0'
9 end
10

```

ภาพที่ 38 แสดงตัวอย่างไฟล์ Podfile ที่บรรจุส่วนเสริม

จากนั้นตั้งชื่อไฟล์ว่า Podfile และทำการบันทึกในโฟลเดอร์ของเครื่องมือ

- 3) จากนั้นใช้คำสั่งต่อไปนี้บนเทอร์มินอล
\$ pod install

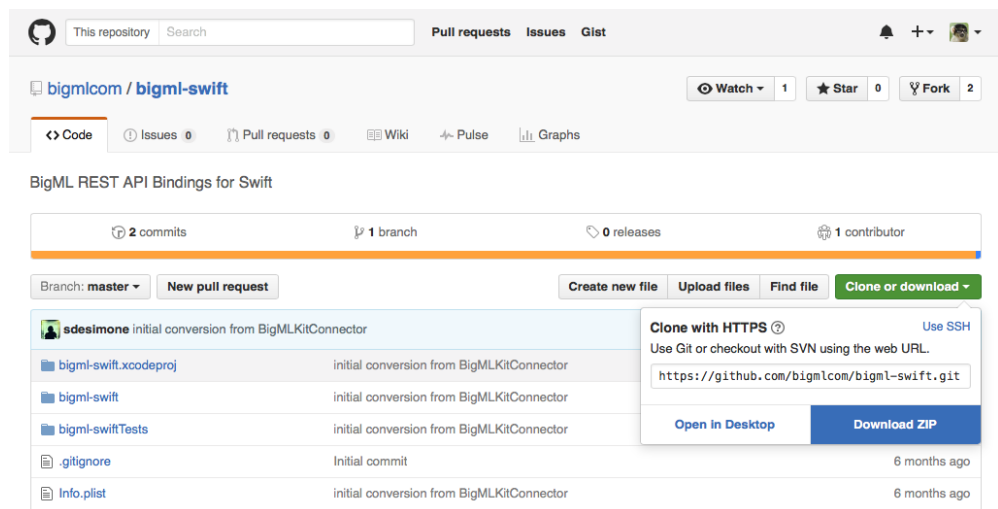
ภาคผนวก ข.

การติดตั้งส่วนต่อประสานเว็บบริการ (bigml-swift)

ในการติดต่อการทำงานกับเว็บบริการ เครื่องมือจำเป็นที่จะต้องมีการนำเข้าสู่ส่วนต่อประสานกับเว็บบริการ ซึ่งมีวิธีการติดตั้งดังต่อไปนี้

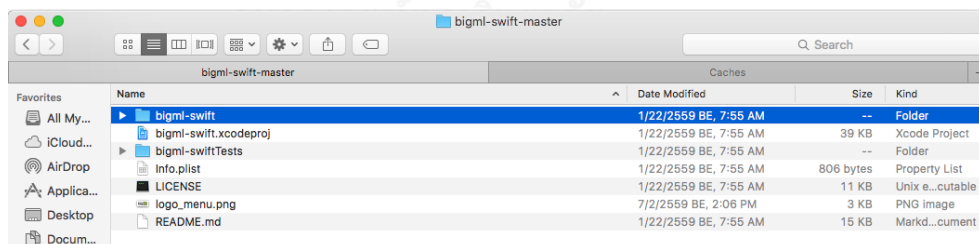
- 1) ทำการดาวน์โหลดส่วนต่อประสานในภาษา Swift จากที่อยู่เว็บ

<https://github.com/bigmlcom/bigml-swift>



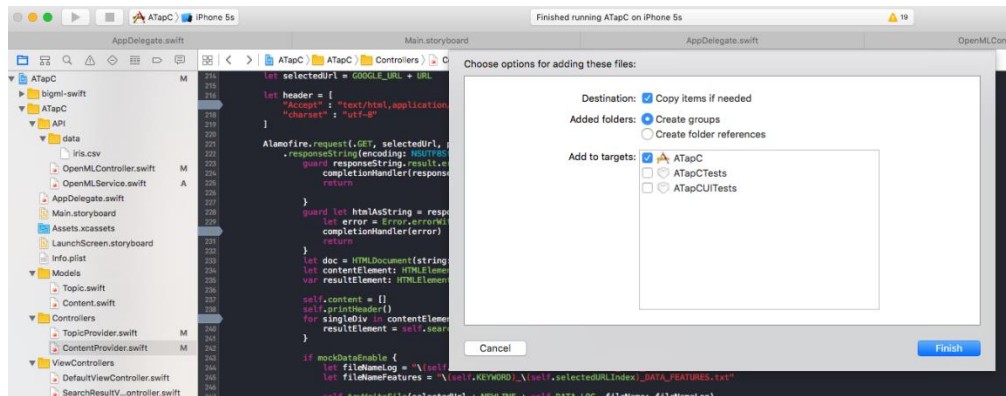
ภาพที่ 39 วิธีการดาวน์โหลดส่วนต่อประสาน

- 2) จากนั้นทำการเปิดเครื่องมือบนโปรแกรม XCode
- 3) แลกไฟล์ที่ได้จากการดาวน์โหลดในข้อที่ 1 ด้วยโปรแกรม Zip



ภาพที่ 40 แสดงโครงสร้างไฟล์เดออร์หลังจากทำการแตกไฟล์

- 4) ทำการลากไฟล์เดออร์ที่มีชื่อว่า bigml-swift ลงในโปรแกรมของเครื่องมือจากนั้น คลิกปุ่ม Finish เพื่อทำการนำเข้าสู่ส่วนต่อประสานสู่ตัวเครื่องมือ



ภาพที่ 41 แสดงหน้าจอในการนำเข้าสู่ส่วนต่อประสาน



ประวัติผู้เขียนวิทยานิพนธ์

นาย กฤษณ์ บรรณะชัยศิริสุข เกิดเมื่อวันที่ 10 มีนาคม พ.ศ. 2532 ที่จังหวัด กรุงเทพมหานคร สำเร็จการศึกษาปริญญาตรีหลักสูตรวิศวกรรมศาสตรบัณฑิต (วศ.บ.) ภาควิชา วิศวกรรมคอมพิวเตอร์ (หลักสูตรนานาชาติ) คณะวิศวกรรมศาสตร์ มหาวิทยาลัยพระจอมเกล้า ธนบุรี ในปีการศึกษา 2554 และเข้าศึกษาต่อในหลักสูตรวิทยาศาสตรมหาบัณฑิต สาขา วิทยาศาสตร์คอมพิวเตอร์ ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์ มหาวิทยาลัย ในปีการศึกษา 2556

