

บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

จากการศึกษาขั้นตอนกระบวนการรู้จำตัวอักษรที่ผ่านมาในงานวิจัยที่เกี่ยวข้องนั้นได้มีการศึกษาวิจัยกันมาอย่างต่อเนื่องและกว้างขวางเพื่อเพิ่มประสิทธิภาพความสามารถในการรู้จำตัวอักษรให้มีความถูกต้องมากขึ้นอาทิเช่น การรู้จำตัวอักษรโดยอาศัยโครงข่ายประสาทเทียมรูปแบบต่าง ๆ ได้แก่ Back-Propagation, BAM และ CPN หรือจะเป็นการรู้จำโดยอาศัยการใช้ลักษณะเด่นของตัวอักษรซึ่งอาจจะมีการนำเอาทฤษฎีฟิชซ์ลอกิหรือทฤษฎีอื่น ๆ เข้ามาช่วยในกระบวนการรู้จำ ทั้งนี้ความสามารถในการรู้จำนั้นไม่ได้ขึ้นกับวิธีที่นำมาใช้ในขั้นตอนของการรู้จำตัวอักษรเท่านั้น แต่ยังต้องอาศัยกระบวนการก่อนการรู้จำที่ดีเพื่อให้ได้ภาพตัวอักษรที่มีคุณภาพดีด้วย ด้วยเหตุนี้จึงทำให้มีข้อดีและข้อเสียแตกต่างกันไป โดยเฉพาะถ้าภาพตัวอักษรก่อนการรู้จำมีคุณภาพที่ไม่ดีก็จะส่งผลให้ประสิทธิภาพในการรู้จำภาพตัวอักษรนั้นไม่มีความถูกต้องเท่าที่ควร[1] ดังนั้นในการรู้จำตัวอักษรก็จะขึ้นอยู่กับกระบวนการก่อนการรู้จำที่จะจัดการข้อมูลภาพตัวอักษรให้อยู่ในรูปแบบที่เหมาะสมก่อนจะเข้าสู่กระบวนการรู้จำ เช่นการปรับให้คุณภาพของข้อมูลภาพตัวอักษร ก่อนการรู้จำมีคุณภาพที่ดีโดยจะต้องไม่มีสัญญาณรบกวน และรวมถึงปัญหาต่าง ๆ ที่พบเจอในงานวิจัยที่ผ่านมา เช่น การแก้ไขตัวอักษรที่ติดกัน [2] [3] [4] การลดสัญญาณรบกวน [4] การทำความสะอาดตัวอักษรให้เรียบ [5] การเชื่อมต่อลายเส้นที่ขาดหายไปของภาพอักษรตัวพิมพ์ภาษาไทยในแนวดิ่ง [7] เหล่านี้ช่วยให้กระบวนการก่อนการรู้จำได้ข้อมูลภาพตัวอักษรที่ดีขึ้น จากนั้นจะพิจารณากระบวนการรู้จำที่มีการออกแบบอัลกอริทึมสำหรับการรู้จำภาษาต่าง ๆ โดยแต่ละรูปแบบภาพตัวอักษรของแต่ละภาษาจะมีคุณลักษณะโครงสร้างที่แตกต่างกันไป สำหรับงานวิจัยที่ผ่านมานั้น จะเป็นการศึกษาในการรู้จำตัวอักษรของแต่ละภาษาโดยเฉพาะ ตัวอย่างเช่นการรู้จำตัวอักษรภาษาอังกฤษ การรู้จำตัวอักษรภาษาจีน การรู้จำตัวอักษรภาษาฝรั่งเศส หรือแม้กระทั่งการรู้จำตัวอักษรภาษาไทย

สำหรับรูปแบบโครงสร้างของตัวอักษรภาษาไทยนั้นปกติแล้วจะประกอบไปด้วยเส้นตัวอักษรที่เป็นเส้นตรง เส้นโค้ง เส้นซิกแซก หรือลูป โดยส่วนมากแล้วตัวอักษรภาษาไทยจะมีจุดเริ่มต้นเป็นลูป หรือที่เราเรียกว่าส่วนที่เป็นหัวตัวอักษร แต่ในปัจจุบันมีการพัฒนารูปแบบตัวอักษรไทยที่เปลี่ยนไปจากเดิมคือที่ส่วนหัวของตัวอักษรไม่ได้เป็นลูปเรียกว่ารูปแบบตัวอักษรที่ไม่มีหัว ตามตารางที่ 1.1 แสดงการเปรียบเทียบลักษณะของตัวอักษรภาษาไทยในรูปแบบที่มีหัว

และไม่มีหัว ดังนั้นมีการพิจารณาลักษณะเด่นของโครงสร้างตัวอักษรภาษาไทย สิ่งแรกที่จะใช้พิจารณาคือลักษณะของตัวอักษรภาษาไทยในรูปแบบที่มีหัว ซึ่งงานวิจัยที่เกี่ยวกับการรู้จำตัวอักษรภาษาไทยส่วนมากนั้นไม่ได้นิยมลักษณะของตัวอักษรภาษาไทยในรูปแบบที่ไม่มีหัว ซึ่งสามารถพบเจอได้ตามเอกสารต่าง ๆ เช่นหนังสือพิมพ์ วารสาร สื่อแผ่นพับ โฆษณาต่าง ๆ ที่เป็นตัวอักษรภาษาไทยที่มีรูปแบบการเขียนหรือพิมพ์ออกแบบแล้วเป็นรูปแบบที่ไม่มีหัวเป็นส่วนใหญ่

ตัวอย่างข้อมูลภาพตัวอักษรตามรูปที่ 1.1 เมื่อผ่านการทดสอบด้วยโปรแกรมอ่านไทย (ArnThai 2.5 Lite) พบว่าผลที่ได้ออกมา มีความถูกต้องเท่ากัน 34.85% และคงผลลัพธ์ที่ได้จากการรู้จำตามรูปที่ 1.2 ซึ่งแสดงให้เห็นว่าผลของการรู้จำที่ได้ออกมาซึ่งคงมีปัญหาในเรื่องของตัวอักษรพิมพ์ภาษาไทยในรูปแบบที่ไม่มีหัว

ในงานวิจัยครั้งนี้จึงเน้นที่ก่อรุ่นข้อมูลตัวอย่างลักษณะของตัวอักษรภาษาไทยในรูปแบบที่ไม่มีหัวซึ่งจะมีทั้งหมด 39 รูปแบบ แต่ละรูปแบบจะประกอบด้วยพยัญชนะ สาระ วรรณยุกต์ ตัวเลข อารบิก และสัญลักษณ์อื่น ๆ รวมทั้งหมดครูปแบบละ 80 ตัวอักษร นอกจากนี้ยังพิจารณากรุ่นข้อมูลภาพตัวอักษรที่อยู่บนป้ายทะเบียนรถที่เป็นภาษาไทยและตัวเลขอารบิก

แต่ทั้งนี้จะเห็นว่าตั้งแต่อดีตมาจนถึงปัจจุบัน ได้มีนักวิจัยทำการพัฒนาเทคโนโลยีทางด้านการรู้จำตัวอักษรกันมากขึ้น แต่ในการพัฒนาจะมุ่งเน้นก่อรุ่นทดลองตามภาษาที่ใช้ของแต่ละประเทศ ด้วยเหตุนี้ การรู้จำของแต่ละภาษา ก็มีวิธีในการรู้จำที่ต่างกันออก ไปตามโครงสร้างของแต่ละภาษา เนื่องมาจากตัวอักษรของแต่ละภาษาที่มีรูปแบบต่างกันออก ไป และสำหรับงานวิจัยที่เกี่ยวกับการรู้จำภาษาไทยนั้น จะรู้จำในก่อรุ่นของตัวอักษรไม่มีหัว แต่ส่วนมากในการทำเลื่อนโยบายฯ หรือ หนังสือต่าง ๆ จะประกอบไปด้วยตัวอักษรภาษาไทยไม่มีหัว ซึ่งปัจจุบันนี้พบอยู่อย่างแพร่หลาย

ดังนั้นงานวิจัยครั้งนี้ได้เลือกเห็นถึงความสำคัญของการรู้จำตัวอักษรตัวพิมพ์ภาษาไทยไม่มีหัว เพื่อเพิ่มประสิทธิภาพในการรู้จำตัวอักษรภาษาไทยให้ครอบคลุมไปถึงตัวอักษรภาษาไทยไม่มีหัว ให้สามารถได้ผลลัพธ์เอกสารสื่อต่าง ๆ ออกแบบในรูปแบบที่ยอมรับ ได้ และมีประสิทธิภาพยิ่งขึ้น

1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา

1.2.1 เพื่อศึกษารากฐานพิเศษของตัวพิมพ์อักษรไทยในก่อรุ่นตัวอักษรไม่มีหัว ที่จะนำมาใช้ในวิธีการรู้จำ โดยศึกษาเปรียบเทียบจากวิธีต่าง ๆ ที่ผ่านการวิจัยมาแล้ว

1.2.2 เพื่อศึกษาวิธีการที่เหมาะสมของการวิเคราะห์หาความสัมพันธ์ของโครงสร้างทางภาษาพามาใช้กับการรู้จำตัวพิมพ์อักษรไทยในก่อรุ่นตัวอักษรไม่มีหัว

1.2.3 เพื่อศึกษาวิธีการที่เหมาะสมของการนำโครงข่ายประสาทเทียมมาใช้กับการรู้จำตัวพิมพ์อักษรไทย

1.2.4 เพื่อเพิ่มประสิทธิภาพในระบบการรู้จำของตัวอักษรภาษาไทยครอบคลุมในรูปแบบตัวอักษรไม่มีหัว

ตารางที่ 1.1 การเปรียบเทียบลักษณะของตัวอักษรภาษาไทยในรูปแบบที่มีหัวและไม่มีหัว

ตัวอักษรที่	ตัวอักษรภาษาไทยรูปแบบที่มีหัว	ตัวอักษรภาษาไทยรูปแบบที่ไม่มีหัว
1. กอไก'	ก	ก
2. ขอไข่	ຂ	ຂ
3. ขอขาวด	ງ	ຂ
4. គុគ្រាយ	គ	ទ
5. កែចន	គ	ទ
6. មូរជំង	ឃ	ឃ
7. សង្ស	ង	ឃ
8. ខខាន	ច	ច
9. ននីង	ឈ	ឈ
10. ខខោង	ឃ	ឃ
11. ធមូទិច	ឃ	ឃ
12. ណុកភេះលោន	ឃ	ឃ
13. សូអូអូហូូុង	ឃ	ឃ
14. នូបត្បូក	ឃ	ឃ
15. នូខ្សាតា	ឃ	ឃ
16. នូខ្សាន	ឃី	ឃី
17. ទុននុពុទា	ឃ	ក
18. ធមូសូត្រោះ	ឃ	ឃ
19. ណុលេនវ	ឃ	ឃ
20. គុគៀក	ទ	ទ
21. គុគោះ	ទ	ទ

ตารางที่ 1.1 (ต่อ)

ตัวอักษรที่	ตัวอักษรภาษาไทยรูปแบบที่มีหัว	ตัวอักษรภาษาไทยรูปแบบที่ไม่มีหัว
22. ถอถุง	ถ	ດ
23. ทอทหาร	ທ	ກ
24. ช้อซัง	ຊ	ຊ
25. นอหนู	ນ	ນ
26. บอใบไม้	ບ	ບ
27. ปอปคลา	ປ	ປ
28. ຜອຜິ້ງ	ຜ	ຜ
29. ຜອຝາ	ຜ	ຜ
30. ພອພານ	ພ	ພ
31. ພອຟິນ	ພ	ພ
32. ກອສໍາເກາ	ກ	ກ
33. ມອມໍາ	ມ	ມ
34. ຍອຍັກຍົງ	ຍ	ຍ
35. ຮອເວືອ	ຮ	ຮ
36. ລອຄິງ	ລ	ລ
37. ວອແຫວນ	ວ	ວ
38. ສອສາລາ	ສ	ສ
39. ພອນອຸນີ່ງ	ໝ	ໝ
40. ສອເຕື່ອ	ສ	ສ
41. ໂອທຶນ	ທ	ທ
42. ພອຈຸພາ	ພ	ຜ
43. ອອອ່າງ	ອ	ວ
44. ຂອນກສູກ	ຂ	ຮ

ກລຸ່ມກາພຕ້ວອັກນຈຣ Jasmine

**ກບຂດຕະງຈອໜະນີນິມີເຈົ້າກົມບົດຕະກຫ
ບບປັພິພຳກມຍຣລວຕະຫສຫພອຍ**

๑๒๓๔๕໬ຕັດໝາຍ, ພະ+ສະການນະພາບໂໄກ

ກລຸ່ມກາພຕ້ວອັກນຈຣ Kodchiang

**ກບຂດຕະງຈອໜະນີນິມີເຈົ້າກົມບົດຕະກຫບປັພ
ພຳກມຍຣລວຕະຫສຫພອຍສານໂໄກ**

ກລຸ່ມກາພຕ້ວອັກນຈຣ Lily

**ກບຂດຕະງຈອໜະນີນິມີເຈົ້າກົມບົດຕະກຫ
ຮນບປັພວິທີກມຍຣລວຕະຫສຫພວວະ**

ກລຸ່ມກາພຕ້ວອັກນຈຣ Fixedsys

**ກບຂດຕະງຈອໜະນີນິມີເຈົ້າກົມບົດຕະກຫ
ຮນບປັພວິທີກມຍຣລວຕະຫສຫພວວະ**

ຮູບທີ 1.1 ຕ້າວຍ່າງຮູບແບບຂໍ້ມູນກາພຕ້ວອັກນຈຣທີ່ໃຊ້ໃນການທົດລອງ

ກລຸ່ມກາພຕ້ວອັກນຈຣ Jasmine

ກບຂດຕະງຈອໜະນີນິມີເຈົ້າກົມບົດຕະກຫບປັພພຳກມຍຣລວຕະຫສຫພວວະ

ກລຸ່ມກາພຕ້ວອັກນຈຣ Kodchiang

ກບຂດຕະງຈອໜະນີນິມີເຈົ້າກົມບົດຕະກຫບປັພພຳກມຍຣລວຕະຫສຫພວວະ

ກລຸ່ມກາພຕ້ວອັກນຈຣ Lily

ນີ້ແມ່ນບັນດາຂອງຂໍ້ມູນຂອງຂໍ້ມູນຂອງຂໍ້ມູນຂອງຂໍ້ມູນຂອງຂໍ້ມູນຂອງຂໍ້ມູນ

ກລຸ່ມກາພຕ້ວອັກນຈຣ Fixedsys

ນີ້ແມ່ນບັນດາຂອງຂໍ້ມູນຂອງຂໍ້ມູນຂອງຂໍ້ມູນຂອງຂໍ້ມູນຂອງຂໍ້ມູນ

ຮູບທີ 1.2 ຕ້າວຍ່າງແສດງພົກງານຮູ້ຈຳເນື້ອຜ່ານໂປຣແກຣມກາຮູ້ຈຳ ArnThai 2.5 Lite

1.3 ทฤษฎีหรือแนวความคิดที่ใช้ในการวิจัย

แนวความคิดที่ใช้ในงานวิจัยนี้ประกอบไปด้วย 2 กระบวนการคือกระบวนการก่อนการรู้จำและกระบวนการรู้จำ ในขั้นตอนของการรู้จำนี้เมื่อรับเอกสารตัวอักษรในเอกสารเข้ามาแล้วจะทำการแบ่งแต่ละบรรทัดออกจากกันด้วยวิธีของการทำ Histogram จากนั้นในแต่ละบรรทัดจะแบ่งออกเป็น 3 ระดับซึ่งจะทำการจัดเรียงกลุ่มแต่ละกลุ่มตามระดับของภาพตัวอักษรนั้น ๆ ได้แก่ กลุ่มตัวอักษรระดับบน กลางและล่างตามลำดับ แล้วใช้วิธี Histogram อีกเช่นกันในการแบ่งภาพตัวอักษรให้ได้ออกมาเป็นกรอบภาพตัวอักษรแต่ละตัว จากนั้นจะผ่านกระบวนการ Thinning กรอบภาพตัวอักษร ขั้นตอนต่อไปคือการแบ่งภาพตัวอักษรที่ Skeleton และวิธีการเป็น 9 ส่วนย่อยที่เท่า ๆ กัน แล้วทำการพิจารณาหาจุดปลายของภาพตัวอักษรและมีการเก็บข้อมูลโดยริบจากจุดปลายที่ 1 ไปจนถึงจุดปลายที่ 2 และหากกรอบส่วนย่อยที่ถูกแบ่งนั้นมีจุดปลายที่สามก็จะอยู่ในกลุ่มของตัวอักษรที่มีจุดปลายมากกว่า 2 จุด ซึ่งค่าที่เก็บนั้นจะมีการแทนรหัสแทนทิศทาง 9 รหัส ตามทิศทางที่ทำมุมของเส้นตัวอักษร โดยที่แต่ละกรอบที่ถูกแบ่งย่อย 9 กรอบจะเก็บค่าที่อยู่ในช่วงกรอบของตนเอง และทำการเรียงเก็บข้อมูลต่อไป กล่าวคือรูปแบบ Input Vector จะเก็บเฉพาะส่วนข้อมูลที่เป็นลายเส้นภาพตัวอักษร จากนั้นแต่ละกรอบที่ถูกแบ่งออกจะส่งต่อเข้าสู่กระบวนการรู้จำโดยอาศัยโครงข่ายประสาทเทียมแบบใหม่ที่พัฒนาขึ้นมาโดยเน้นที่บริเวณส่วนไม่มีหัวของตัวอักษรภาษาไทยมาเป็นหลักการเบื้องต้นของแนวคิดในการพัฒนากระบวนการเรียนรู้ทั้งหมด และพยายามที่จะสร้างกระบวนการรู้จำใหม่ที่มีการปรับเปลี่ยนเฉพาะบางส่วนย่อยเพื่อให้เกิดประสิทธิภาพในการเรียนรู้บริเวณตัวอักษรไม่มีหัวที่เป็นจุดเริ่มต้นของการเรียนรู้ดังนั้น นอกจากระบบทั่วไปนี้ได้นำเอาวิธีการวัดความเหมือนหรือความคล้ายคลึงกันที่อาศัยฟังก์ชันกรอสโครริเลชัน (Cross-correlation) ที่มีประสิทธิภาพและสามารถนำมาเป็นตัววัดค่าความเหมือนของรูปแบบ ส่วนแนวทางในการรู้จำ Fuzzy Adaptive Resonance Theory (Fuzzy ARTMAP) หรือที่เรียกว่าฟิชช์อาร์ทแมปนั้น มีความสามารถในการเรียนรู้สิ่งใหม่ได้พร้อมกับยังคงความรู้จำเก่าไว้ได้ ดังนั้นงานวิจัยนี้จึงใช้แนวทางฟิชช์อาร์ทแมปเป็นพื้นฐานของกระบวนการเรียนรู้และพยายามหาวิธีการปรับเอกสารกรอสโครริเลชันมาใช้ร่วมกับแนวทางฟิชช์อาร์ทแมปให้ได้

1.4 ขอบเขตของการวิจัย

งานวิจัยนี้เน้นการพัฒนากระบวนการรู้จำสำหรับรูปแบบภาพตัวอักษรในเอกสารที่มีขอบเขตต่อไปนี้

1.4.1 ภาพตัวอักษรเอกสารที่ได้มาจากการสแกนเข้าสู่เครื่องคอมพิวเตอร์ ภาพตัวอักษรจะต้องไม่คล้ายกันทั้งภายใน-นอก ตัวอักษร รวมไปถึงไม่มีการซ้อนทับกันของภาพตัวอักษร และภาพลายเส้นของตัวอักษรจะต้องไม่ขาด

1.4.2 ภาพตัวอักษรเอกสารที่ได้มาจากการสแกนเข้าสู่เครื่องคอมพิวเตอร์ จะต้องประกอบด้วยตัวอักษรภาษาไทยไม่มีหัว โดยครอบคลุมตัวอักษรตามตารางที่ 1.2 ต่อไปนี้

ตารางที่ 1.2 ตัวอักษรภาษาไทยฟอนต์ไม่มีหัว

พยัญชนะ 46 ตัว	ก ข ช ດ ຕ ນ ຈ ຈ ະ ນ ິ ງ ີ ື ຶ ດ ຕ ກ ກ ຮ ບ ປ ຜ ົ ພ ພ ກ ມ ຍ ລ ວ ຕ ສ ທ ຜ ພ ອ ກ ກ
สาระ 15 ตัว	ຂ ໂ ກ ແ ່ ໊ ເ ໃ ໄ ແ ້ ໂ ໆ ່ ້
วรรณยุกต์ 7 ตัว	໌ ໂ ແ ່ ໄ ້ ໊
เลขอารบิก 10 ตัว	0 1 2 3 4 5 6 7 8 9
สัญลักษณ์ 2 ตัว	໌ ໍ

1.5 ขั้นตอนของการศึกษา

1.5.1 ศึกษานบทความและผลงานวิจัยต่างๆ ที่มีความเกี่ยวข้องกับงานวิจัยนี้

1.5.2 เก็บข้อมูลตัวอย่างของตัวอักษร พร้อมจัดเก็บลงคอมพิวเตอร์

1.5.3 ศึกษาลักษณะโครงสร้างของตัวอักษรภาษาไทยเพื่อนำไปวิเคราะห์

1.5.4 ออกแบบอัลกอริทึมในการรู้จำโดยใช้ Hierarchical Cross-Correlation Neural Network

1.5.5 เก็บนิปปองโปรแกรมเพื่อวิเคราะห์ลักษณะของภาพตัวอักษร และทำการจัดกลุ่มของตัวอักษรในส่วนของกระบวนการก่อนการรู้จำ

1.5.6 เก็บนิปปองโปรแกรมเพื่อรู้จำภาพตัวอักษร โดยใช้ Hierarchical Cross-Correlation Neural Network

1.5.7 ทดลองรู้จำตัวอักษรกับข้อมูลที่จัดเก็บ

1.5.8 ทดสอบผลจากการรู้จำตัวอักษรว่ามีความถูกต้องเมื่อเทียบกับซอฟต์แวร์ที่ใช้งานจริง ArnThai version 2.5 Lite และ Thai OCR โดยทดสอบข้อมูลก่อนและหลังการรู้จำตัวอักษร

1.5.9 สรุปผลการดำเนินการ และรวมรวมนำจัดทำเอกสารนำเสนอเป็นงานวิจัย

1.6 ประโยชน์ที่คาดว่าจะได้รับ

1.6.1 ทราบถึงบทความงามนิจย์ต่าง ๆ ที่เกี่ยวข้องกับงานนิจย์ครั้งนี้

1.6.2 ทราบถึงลักษณะเด่นของตัวพิมพ์อักษรภาษาไทยแบบฟอนต์ไม่มีส่วนหัวของตัวอักษร

1.6.3 เป็นพื้นฐานในการรู้จำตัวอักษรภาษาไทยที่สมบูรณ์แบบ เพื่อใช้ในการพัฒนาต่อไปในอนาคตให้เหมาะสมกับตัวอักษรภาษาไทยทุกรูปแบบ