

การทดสอบการรู้จำตัวพิมพ์อักษรภาษาไทย

ในบทนี้จะกล่าวถึงการทดสอบผลการรู้จำ โดยใช้กฎที่ PROGOL เรียนรู้ได้จากที่กล่าวไว้ในบทที่ 3 ซึ่งรายละเอียดในบทนี้ประกอบด้วย วิธีการทดสอบ ผลการทดสอบ และแสดงให้เห็นถึง ปัญหา และข้อจำกัด ในส่วนท้ายของบท

วิธีการทดสอบ

หลังจากการเรียนรู้ได้นิยามของตัวอักษรในรูปของอนุประโยคแล้ว อนุประโยคดังกล่าวจะถูกเก็บในแฟ้มข้อมูล "rule.dat" ในรูปแบบที่กำหนด โดยโปรแกรมที่ใช้ทดสอบการรู้จำจะอ่านนิยามของตัวอักษรจากแฟ้มข้อมูลนี้ โปรแกรมที่ใช้ทดสอบการรู้จำพัฒนาบน Borland C++ รุ่น 3.1 การทำงานของโปรแกรมประกอบด้วย การรับแฟ้มข้อมูลภาพของตัวอักษรที่ใช้ในการทดสอบ แล้วนำมาผ่านการประมวลผลขั้นต้น และ เปรียบเทียบลักษณะของตัวอักษรกับนิยามของตัวอักษรที่ได้จากการเรียนรู้ในแฟ้มข้อมูล "rule.dat" ถ้าลักษณะของตัวอักษรตรงกับนิยามใด ก็จะได้ผลการรู้จำเป็นตัวอักษรนั้น แต่ถ้าลักษณะของตัวอักษรไม่ตรงกับนิยามของตัวอักษรใดเลย โปรแกรมจะเลือกนิยามของตัวอักษรที่ใกล้เคียงที่สุดให้เป็นผลการรู้จำ โดยมีวิธีการเลือก ดังนี้

1. เลือกนิยามของตัวอักษรที่มีระดับของตัวอักษรอยู่ในระดับเดียวกันเท่านั้น
2. เลือกนิยามที่มีจำนวนเพรดิเคต ซึ่งไม่ตรงกับลักษณะของตัวอักษรที่ใช้ทดสอบ มีจำนวนน้อยที่สุด
3. ในกรณีที่ได้ผลจากข้อ 1 มากกว่า 1 นิยาม ให้เลือกนิยามที่มีจำนวนเพรดิเคต ซึ่งตรงกับลักษณะของตัวอักษรที่ใช้ทดสอบ มีจำนวนมากที่สุด

ตัวอย่าง สมมติให้มีกฎในการรู้จำตัวอักษร 4 ตัวอักษร คือ fan(ฟ) sampoa(ก) mar(ม) yag(ย) และ ข้อมูลตัวอักษรที่นำมารู้จำไม่ตรงกับนิยามใดเลย จึงต้องเลือกนิยามที่ใกล้เคียงที่สุด โดยที่ ตัวอักษรที่นำมารู้จำมีระดับของตัวอักษรอยู่ในระดับที่ 3

เครื่องหมาย ✓ บอกถึงเพรดิเคตที่ตรงกับลักษณะของตัวอักษรที่นำมารู้จำ

เครื่องหมาย ✗ บอกถึงเพรดิเคตที่ไม่ตรงกับลักษณะของตัวอักษรที่นำมารู้จำ

fan(2,A,B,C,D,E) :- enptprim(B,1), enptprim(B,6).

x

✓

✓

sampo(3,A,B,C,D,E) :- headzone(B,3), headprim_0009(B), memberzone(B,1,4).

✓

✓

✓

x

mar(3,A,B,C,D,E) :- headzone(B,2), headprim(B,12), enptzone(B,4), enptprim(B,1).

✓

x

x

✓

✓

yag(3,A,B,C,D,E) :- havemember(B,4,2,2), cntstzone1(B,0).

✓

x

✓

วิธีการเลือก คือ

1. เลือกเฉพาะนิยามที่มีระดับของตัวอักษรถูกต้อง ได้แก่ sampo, mar และ yag
2. เลือกนิยามที่มีเพรดิเคต x น้อยที่สุด ได้แก่ sampo และ yag ซึ่งมีจำนวน 1 เพรดิเคต
3. เลือกนิยามที่มีเพรดิเคต ✓ มากที่สุด ได้แก่ sampo มีจำนวน 3 เพรดิเคต ในขณะที่ yag มีจำนวน 2 เพรดิเคต ดังนั้น จึงได้ผลการรู้จำ คือ ตัวอักษร 'ภ' (sampo)

ผลการทดสอบ

จากการเรียนรู้ซึ่งแบ่งเป็น 2 กลุ่มดังที่กล่าวมาแล้ว ผลการเรียนรู้ในแต่ละกลุ่มจะนำมาใช้ในการทดสอบตัวอักษรที่แตกต่างกัน ดังนี้

1. การรู้จำตัวอักษรในรูปแบบที่ไม่เคยผ่านการเรียนรู้มาก่อน ข้อมูลที่ใช้ในการเรียนรู้ (กลุ่มที่ 1) คือ ตัวพิมพ์อักษรภาษาไทย รูปแบบ EUCROSIA ขนาด 48, 36, 32, 28, 24, 22 และ 20

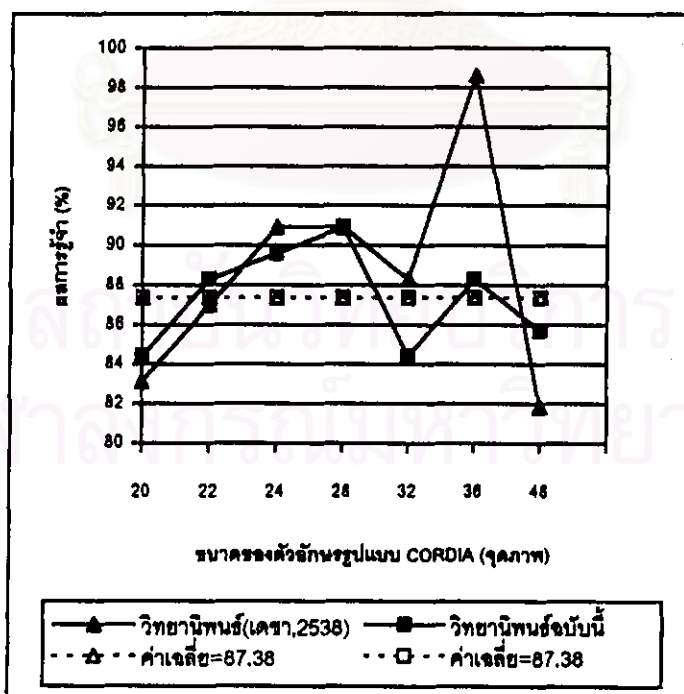
ข้อมูลที่ใช้ในการทดสอบ คือ ตัวพิมพ์อักษรภาษาไทย 77 ตัวอักษร รูปแบบ CORDIA ขนาด 48, 36, 32, 28, 24, 22 และ 20 รวม 539 ตัวอักษร ที่มาของข้อมูลได้จากการป้อนตัวอักษรทั้งหมดที่มีรูปแบบและขนาดตามที่ต้องการ โดยใช้โปรแกรม Microsoft Word 97 แล้วพิมพ์ออกมาด้วยเครื่องพิมพ์เลเซอร์ ที่ความละเอียดขนาด 300 dpi นำเอกสารที่ได้มาทำการอ่านกลับเป็นแฟ้มข้อมูลภาพด้วยเครื่องกวาดตรวจ ที่ความละเอียดขนาด 300 dpi จากนั้นใช้โปรแกรม Paintbrush ของ Microsoft Windows 95 ในการคัดตัวอักษรเป็นตัวอักษรเดี่ยวเก็บลงแฟ้มข้อมูลภาพ 1 แฟ้มต่อ 1 ตัวอักษร แล้วนำไปทดสอบการรู้จำ

จากการทดสอบในกลุ่มที่ 1 ทำให้รู้ถึงความสามารถในการเรียนรู้โดยใช้ระบบไฮแอลที ว่าสามารถรู้จำตัวอักษรในรูปแบบที่ไม่เคยผ่านการเรียนรู้มาก่อน ได้ความถูกต้องเฉลี่ย 87.38% ดังแสดงในตารางที่ 4.1 และเปรียบเทียบผลการรู้จำนี้กับงานวิจัยของ เชา รัตนาราช [2] ซึ่งใช้เทคนิคแบบพีซีโลจิก และวิธีซินแทกติก โดยใช้ข้อมูลทดสอบชุดเดียวกัน กราฟแสดงการเปรียบเทียบในรูปแบบที่ 4.1 จากกราฟพบว่า ค่าความถูกต้องเฉลี่ยของทั้งสองงานวิจัยมีค่าเท่ากัน ซึ่งอาจกล่าวได้ว่า เทคนิคการเรียนรู้โดยใช้

ระบบไอแอลที และ การใช้เทคนิคแบบพีซีซีไอจิกและวิธีซินแทกติก มีความสามารถใกล้เคียงกัน ในการรู้จำตัวอักษรในรูปแบบที่ไม่เคยเรียนรู้มาก่อน

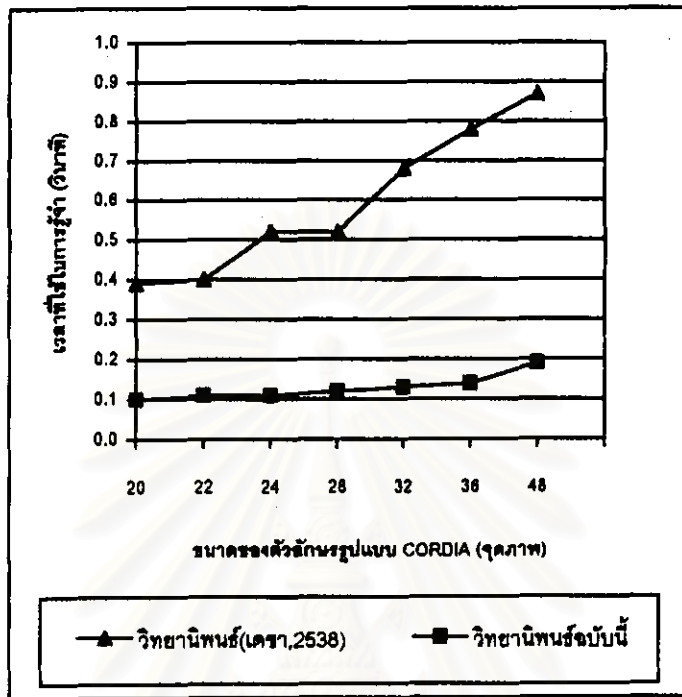
ตารางที่ 4.1 ผลการทดสอบการรู้จำตัวอักษรในรูปแบบที่ไม่เคยผ่านการเรียนรู้มาก่อน

รูปแบบ CORDIA ขนาด	จำนวน	ผลการรู้จำโดยเปรียบเทียบกับกฎที่ได้				กรณีที่ไม่รู้จำไม่ได้ เลือกกฎที่ใกล้เคียงที่สุด		
		รู้จำได้	รู้จำ ไม่ได้	รู้จำผิด	ความถูกต้อง (%)	รู้จำได้	รู้จำผิด	ความถูกต้อง (%)
20	77	64	4	9	83.12	65	12	84.42
22	77	66	4	7	85.71	68	9	88.31
24	77	66	7	4	85.71	69	8	89.61
28	77	68	4	5	88.31	70	7	90.91
32	77	63	6	8	81.82	65	12	84.42
36	77	67	3	7	87.01	68	9	88.31
48	77	61	7	9	79.22	66	11	85.71
รวม	539	455	35	49	84.42	471	68	87.38



รูปที่ 4.1 กราฟเปรียบเทียบเปอร์เซ็นต์ความถูกต้อง การรู้จำตัวอักษรรูปแบบที่ไม่เคยเรียนรู้มาก่อน

จากการทดสอบการรู้จำตัวอักษรในกลุ่มที่ 1 ของงานวิจัยนี้ ใช้เวลาโดยเฉลี่ย 0.13 วินาทีในการรู้จำ 1 ตัวอักษร และการรู้จำโดยใช้งานวิจัยของ เคชา รัตนธาร [2] ใช้เวลาโดยเฉลี่ย 0.59 วินาทีในการรู้จำ 1 ตัวอักษร ดังแสดงกราฟเปรียบเทียบเวลาที่ใช้ในการรู้จำในรูปที่ 4.2



รูปที่ 4.2 กราฟเปรียบเทียบเวลาที่ใช้ในการรู้จำตัวอักษร ในรูปแบบที่ไม่เคยเรียนรู้มาก่อน

2. การรู้จำตัวอักษรที่มีสัญญาณรบกวนปะปน ข้อมูลที่ใช้ในการเรียนรู้ (กลุ่มที่ 2) คือ ตัวพิมพ์อักษรภาษาไทย รูปแบบ CORDIA และ EUCROSIA ขนาด 48, 36, 32, 28, 24, 22 และ 20

ข้อมูลที่ใช้ในการทดสอบ คือ ตัวพิมพ์อักษรภาษาไทย 77 ตัวอักษร รูปแบบ CORDIA และ EUCROSIA ขนาด 48, 36, 32, 28, 24, 22 และ 20 รวม 1078 ตัวอักษร ที่มาของข้อมูลได้จากการป้อนตัวอักษรทั้งหมดที่มีรูปแบบและขนาดตามที่ต้องการ โดยใช้โปรแกรม Microsoft Word 97 แล้วพิมพ์ออกมาด้วยเครื่องพิมพ์เลเซอร์ ที่ความละเอียดขนาด 300 dpi นำเอกสารที่ได้ไปถ่ายเอกสาร โดยถ่ายเอกสาร 2 ครั้งในระดับความเข้มที่ต่างกัน แล้วจึงอ่านกลับเป็นแฟ้มข้อมูลภาพด้วยเครื่องกวาดตรวจ ที่ความละเอียดขนาด 300 dpi จากนั้นใช้โปรแกรม Paintbrush ของ Microsoft Windows 95 ในการตัดตัวอักษรเป็นตัวอักษรเดี่ยวเก็บลงแฟ้มข้อมูลภาพ 1 แฟ้มต่อ 1 ตัวอักษร แล้วนำไปทดสอบการรู้จำ จากการถ่ายเอกสาร 2 ครั้งกับตัวอักษร 1078 ตัว ทำให้ได้ตัวอักษรที่ใช้ทดสอบรวม 2156 ตัวอักษร

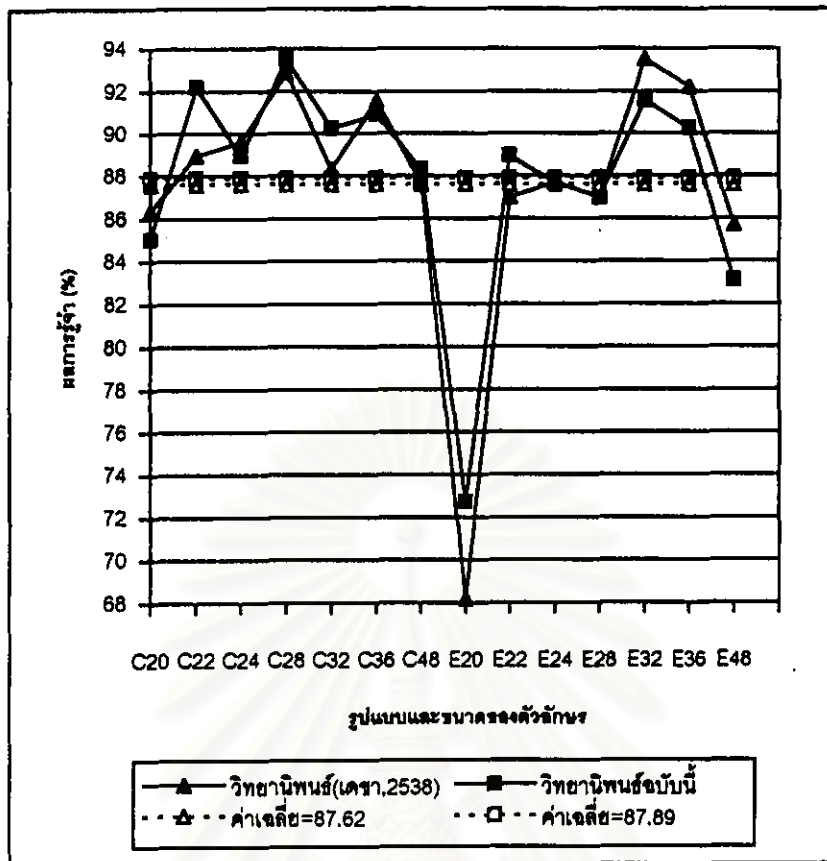
จากการทดสอบในกลุ่มที่ 2 ทำให้รู้ถึงความสามารถในการเรียนรู้ของระบบไอแอลพี ว่า สามารถรู้จำตัวอักษรที่มีสัญญาณรบกวนปะปนอยู่ ซึ่งเกิดจากการถ่ายเอกสาร ได้ความถูกต้องเฉลี่ย 87.89% ดังแสดงในตารางที่ 4.2 และเปรียบเทียบผลการรู้จำนี้กับงานวิจัยของ เศษ รัตนาร [2] ซึ่งใช้เทคนิคแบบฟัชซีโลจิก และวิธีซินแทกติก โดยใช้ข้อมูลทดสอบชุดเดียวกัน กราฟแสดงการเปรียบเทียบในรูปที่ 4.3 จากกราฟพบว่า เทคนิคการเรียนรู้โดยใช้ระบบไอแอลพี มีค่าความถูกต้องเฉลี่ยสูงกว่า การใช้เทคนิคแบบฟัชซีโลจิกและวิธีซินแทกติก ซึ่งมีความถูกต้องเฉลี่ย 87.62% ซึ่งอาจกล่าวได้ว่า การเรียนรู้โดยใช้ระบบไอแอลพีมีความสามารถในการรู้จำตัวอักษรที่มีสัญญาณรบกวนได้ดีกว่า

ตารางที่ 4.2 ผลการทดสอบการรู้จำตัวอักษรที่มีสัญญาณรบกวนปะปน

รูปแบบและขนาด	จำนวน	ผลการรู้จำโดยเปรียบเทียบกับกฎที่ได้				กรณีที่ยังไม่ได้เลือกกฎที่ใกล้เคียงที่สุด		
		รู้จำได้	รู้จำไม่ได้	รู้จำผิด	ความถูกต้อง (%)	รู้จำได้	รู้จำผิด	ความถูกต้อง (%)
C 20	154	124	18	12	80.52	131	23	85.06
C 22	154	138	9	7	89.61	142	12	92.21
C 24	154	134	11	9	87.01	137	17	88.96
C 28	154	136	11	7	88.31	144	10	93.51
C 32	154	136	10	8	88.31	139	15	90.26
C 36	154	138	8	8	89.61	140	14	90.91
C 48	154	131	15	8	85.06	136	18	88.31
รวม (C)	1078	937	82	59	86.92	969	109	89.89
E 20	154	105	29	20	68.18	112	42	72.73
E 22	154	136	7	11	88.31	137	17	88.96
E 24	154	129	16	9	83.77	135	19	87.66
E 28	154	129	18	7	83.77	134	20	87.01
E 32	154	138	13	3	89.61	141	13	91.56
E 36	154	136	15	3	88.31	139	15	90.26
E 48	154	122	21	11	79.22	128	26	83.12
รวม (E)	1078	895	119	64	83.02	926	152	85.90
รวม	2156	1832	201	123	84.97	1895	261	87.89

หมายเหตุ C แทนตัวอักษรรูปแบบ CORDIA

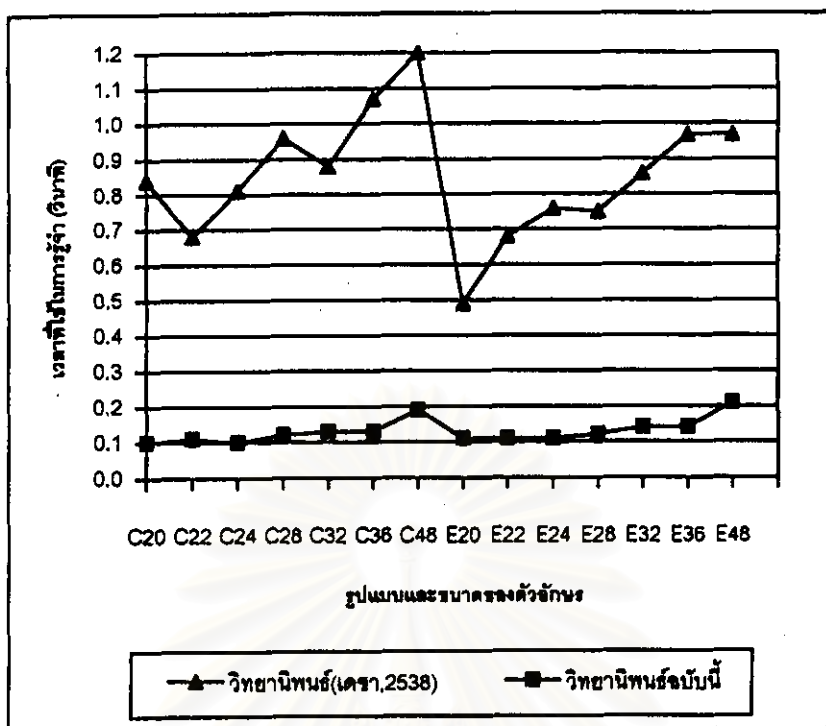
E แทนตัวอักษรรูปแบบ EUCROSIA



รูปที่ 4.3 กราฟเปรียบเทียบเปอร์เซ็นต์ความถูกต้อง การรู้จำตัวอักษรที่มีสัญญาณรบกวนปะปน

จากการทดสอบการรู้จำตัวอักษรในกลุ่มที่ 2 ของงานวิจัยนี้ ใช้เวลาโดยเฉลี่ย 0.13 วินาทีในการรู้จำ 1 ตัวอักษร และการรู้จำโดยใช้งานวิจัยของ เคชา รัตนธาร [2] ใช้เวลาโดยเฉลี่ย 0.85 วินาทีในการรู้จำ 1 ตัวอักษร ดังแสดงกราฟเปรียบเทียบเวลาที่ใช้ในการรู้จำในรูปที่ 4.4

จากการทดสอบผลการรู้จำของตัวอักษรทั้งสองกลุ่ม มีข้อสังเกตประการหนึ่ง คือ งานวิจัยนี้ใช้เวลาในการรู้จำตัวอักษรมากกว่างานวิจัยของ เคชา รัตนธาร เนื่องจาก งานวิจัยของ เคชา รัตนธาร ใช้วิธีการจำแนกตัวอักษร โดยการเปรียบเทียบตัวอักษรที่ต้องการรู้จำกับตัวอักษรต้นแบบ จึงใช้เวลายาวมากกับการอ่านตัวอักษรต้นแบบที่เก็บไว้ในแฟ้มข้อมูลบนงานบันทึกชนิดแข็ง (hard disk) สำหรับงานวิจัยนี้ใช้วิธีการจำแนกตัวอักษร โดยการเปรียบเทียบตัวอักษรที่ต้องการรู้จำกับกฎที่ได้จากการเรียนรู้ ซึ่งกฎดังกล่าวเป็นข้อมูลที่มีขนาดเล็ก จึงใช้เวลาในการอ่านข้อมูลไม่มากนัก รวมทั้งใช้การเปรียบเทียบแบบเงื่อนไขธรรมดา ไม่มีการคำนวณที่ซับซ้อน จึงใช้เวลาในการรู้จำตัวอักษรมากกว่า



รูปที่ 4.4 กราฟเปรียบเทียบเวลาที่ใช้ในการรู้จำตัวอักษรที่มีสัญญาณรบกวนปะปน

ปัญหาและข้อจำกัด

1. ปัญหาเรื่องความละเอียดของจุดภาพ

ความละเอียดของจุดภาพที่ได้จากเครื่องกวาดตรวจ มีผลต่อการรู้จำเป็นอย่างมาก ถ้าใช้ความละเอียด 300 dpi จะมีผลต่อตัวอักษรที่มีขนาดเล็ก ทำให้เส้นของตัวอักษรเลอะ และบางครั้งติดกับเส้นข้างเคียง หรือทำให้ส่วนที่เป็นช่องว่างเต็ม เป็นผลให้รูปโครงร่างของตัวอักษรเปลี่ยนแปลงไป ดังนั้นจึงควรใช้เครื่องกวาดตรวจที่มีความละเอียด 600 dpi จะช่วยแก้ปัญหาได้

2. ปัญหาเรื่องสัญญาณรบกวน จากการถ่ายเอกสาร

ถ้าข้อมูลภาพตัวอักษรได้จากการถ่ายเอกสาร จะทำให้มีสัญญาณรบกวนปะปน บางครั้งมีจุดภาพที่เกิน หรือ บางส่วนของจุดภาพตัวอักษรขาดหายไป เป็นผลให้รูปโครงร่างของตัวอักษรเปลี่ยนแปลงไป ทำให้การรู้จำผิดพลาดได้ ดังนั้น จึงควรให้ความสำคัญกับการถ่ายเอกสาร โดยถ่ายเอกสารให้มีความชัดเจนมากที่สุด จะทำให้ผลการรู้จำมีความถูกต้องมากขึ้น

3. ปัญหาเรื่องลำดับของการกำหนดรูปแบบสัญลักษณ์

การเรียนรู้ตัวอักษรบางตัวของ PROGOL ในงานวิจัยนี้ ได้กฎในการนิยามตัวอักษรมากกว่า 1 กฎ ซึ่งเกิดจากการเรียนรู้ในแต่ละครั้ง มีการจัดลำดับของการกำหนดรูปแบบสัญลักษณ์ที่แตกต่างกัน มีผลทำให้กฎที่ได้จากการเรียนรู้แตกต่างกัน

ตัวอย่าง การเรียนรู้ตัวอักษร 'น' สมมติให้มีเพรดิเคต headzone, headprim และ enptzone

การเรียนรู้ครั้งที่ 1 มีการจัดลำดับการประกาศรูปแบบสัญลักษณ์ ดังนี้

`:- modeb(1,headzone(+sectionlist,#zone))?`

`:- modeb(1,headprim(+sectionlist,#prim))?`

`:- modeb(*,enptzone(+sectionlist,#zone))?`

กฎที่ได้จากการเรียนรู้ คือ `nuu(3,A,B,C,D,E) :- headzone(B,2), headprim(B,12).`

การเรียนรู้ครั้งที่ 2 มีการจัดลำดับการประกาศรูปแบบสัญลักษณ์ต่างจากการเรียนรู้ครั้งแรก ดังนี้

`:- modeb(1,headzone(+sectionlist,#zone))?`

`:- modeb(*,enptzone(+sectionlist,#zone))?`

`:- modeb(1,headprim(+sectionlist,#prim))?`

กฎที่ได้จากการเรียนรู้ คือ `nuu(3,A,B,C,D,E) :- headzone(B,2), enptzone(B,4).`

จากการเรียนรู้ทั้ง 2 ครั้ง แสดงว่า `headzone(B,2), headprim(B,12)` และ `headzone(B,2), enptzone(B,4)` ต่างก็สามารถใช้นิยามตัวอักษร 'น' ซึ่งเป็นตัวอย่างบวกได้ โดยไม่นิยามตัวอักษรอื่นๆ ในตัวอย่างลบ

จากตัวอย่างดังกล่าวข้างต้น พบว่า ลำดับของการประกาศรูปแบบสัญลักษณ์ที่แตกต่างกัน อาจจะมีผลทำให้กฎที่ได้จากการเรียนรู้แตกต่างกัน ซึ่งไม่สามารถทราบได้ว่า กฎใดสามารถใช้นิยามตัวอักษรได้ดีที่สุด ในงานวิจัยนี้ ใช้วิธีการเลือก โดยการนำกฎที่ได้ทั้งหมด มาทดสอบผลการรู้จำกับข้อมูลทดสอบ และเลือกกฎที่ให้ผลการรู้จำดีที่สุด ซึ่งถ้านำไปทดสอบกับข้อมูลชุดอื่น กฎดังกล่าวอาจจะไม่ใช่กฎที่ดีที่สุดก็ได้

จุฬาลงกรณ์มหาวิทยาลัย