

บทที่ 4

การทดสอบการรู้จำตัวอักษรพิมพ์ภาษาไทย

ภาพตัวพิมพ์อักษรภาษาไทย

ภาพข้อมูลของตัวพิมพ์อักษรภาษาไทยที่ใช้ในการทดลอง จะประกอบด้วยตัวพยัญชนะ 44 ตัวอักษร ตัวเลข 10 ตัวอักษร สระ, วรรณยุกต์และตัวอักษรพิเศษ ทั้งหมด 24 ตัวอักษร รวมเป็น 78 ตัวอักษรต่อ 1 ชุดตัวอักษร โดยรูปแบบตัวอักษรที่ใช้ในการทดลองคือ EUCROSIA ขนาด 20, 22, 24, 28, 32, และ 36 จุด ซึ่งภาพข้อมูลตัวอักษรแต่ละชุดของแต่ละขนาดจะประกอบด้วยรูปแบบตัวอักษรปกติ, ตัวเอน และตัวหนา เช่นภาพข้อมูลตัวอักษรขนาด 22 จุด จะประกอบด้วยรูปแบบตัวอักษรปกติ, ตัวเอน และตัวหนา จำนวนทั้งหมด 234 ตัวอักษร ดังนั้นข้อมูลภาพตัวอักษรทั้งหมดที่ใช้ในการทดลอง จะมีด้วยกัน 1392 ภาพตัวอักษร โดยที่ข้อมูลภาพตัวอักษรทั้งหมดจะถูกแบ่งออกเป็น 2 กลุ่ม คือ กลุ่มของชุดฝึก (EUCROSIA 20, 24, 32 point) และกลุ่มของชุดทดสอบ (EUCROSIA 22, 28, 36)

ภาพตัวอักษรทั้งหมดจะได้จากโปรแกรม Microsoft Word ภาษาไทย for Windows 6.0 โดยป้อนตัวอักษรที่ต้องการทั้งหมดโดยใช้โปรแกรมดังกล่าว แล้วพิมพ์ออกมาโดยใช้เครื่องพิมพ์เลเซอร์ที่มีความละเอียด 600 dpi จากนั้นนำภาพตัวอักษรที่ได้จากเครื่องพิมพ์เลเซอร์ไปทำการอ่านกลับมาเป็นไฟล์ของภาพตัวอักษรโดยใช้เครื่องสแกนเนอร์ที่มีความละเอียด 600 dpi หลังจากนั้นจะใช้โปรแกรม Microsoft Windows Paintbrush 3.1 อ่านภาพข้อมูลของตัวอักษรที่สแกนได้จากเครื่องสแกนเนอร์ จากนั้นตัดภาพตัวอักษรเป็นตัวอักษรเดี่ยวๆแล้วจัดเก็บภาพตัวอักษร 1 ภาพตัวอักษรต่อ 1 ไฟล์ข้อมูล โดยจัดเก็บภาพตัวอักษรแต่ละตัวเป็นไฟล์ข้อมูลตระกูลรูปแบบ BMP หรือ PCX เพื่อใช้เป็นภาพข้อมูลสำหรับการฝึกระบบ (training set) และการทดสอบการรู้จำ (test set) ต่อไป



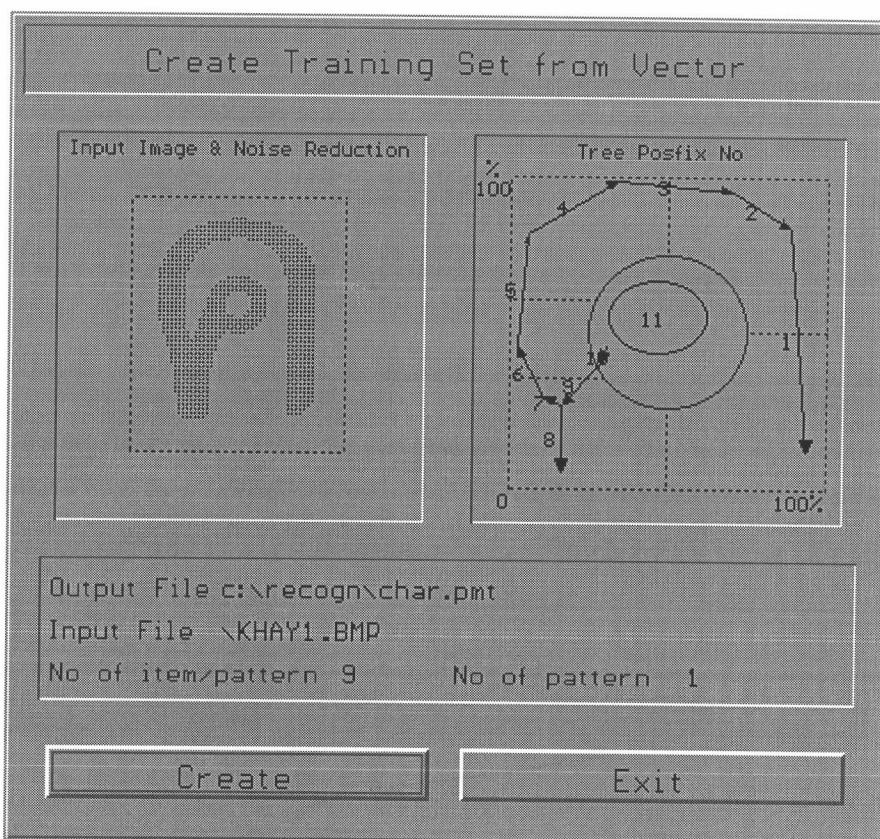
การเรียนรู้ภาพตัวอักษรของระบบนิรอลเน็ตเวิร์ก

ในขั้นตอนแรกก่อนที่ระบบนิรอลเน็ตเวิร์กจะสามารถรู้จำตัวอักษรภาษาไทยได้นั้น จะต้องนำภาพตัวอักษรบางส่วนมาสอนระบบนิรอลเน็ตเวิร์กเสียก่อน โดยภาพข้อมูลที่ใช้เป็นชุดฝึก (training set) คือ ภาพตัวอักษรขนาด 20, 24, 32 จุด ทั้งรูปแบบตัวอักษรปกติ, ตัวเอน และตัวหนา ซึ่งมีข้อมูลภาพตัวอักษรทั้งหมด 690 ภาพตัวอักษร โดยจะแบ่งกลุ่มของตัวอักษรที่มีลักษณะคล้ายคลึงกัน เป็น 11 กลุ่มดังนี้

1. ก ถ ภ ฤ ฎ
2. ฎ ฏ ฒ ฌ ฐ ฌ
3. ค ด ศ ต ค
4. ท ห ฑ ฌ น ม ฆ
5. บ ป ย ข ช ฌ ฌ
6. ฟ ผ ฝ พ ย
7. ล ส จ ฐ ร อ ฮ ว ง
8. ฉ ะ โ ใ ไ ใ ๆ ๆ ๆ ๆ ๆ
9. อ อ อ อ
10. อ อ อ อ อ อ อ
11. ๑ ๒ ๓ ๔ ๕ ๖ ๗ ๘ ๙ ๐

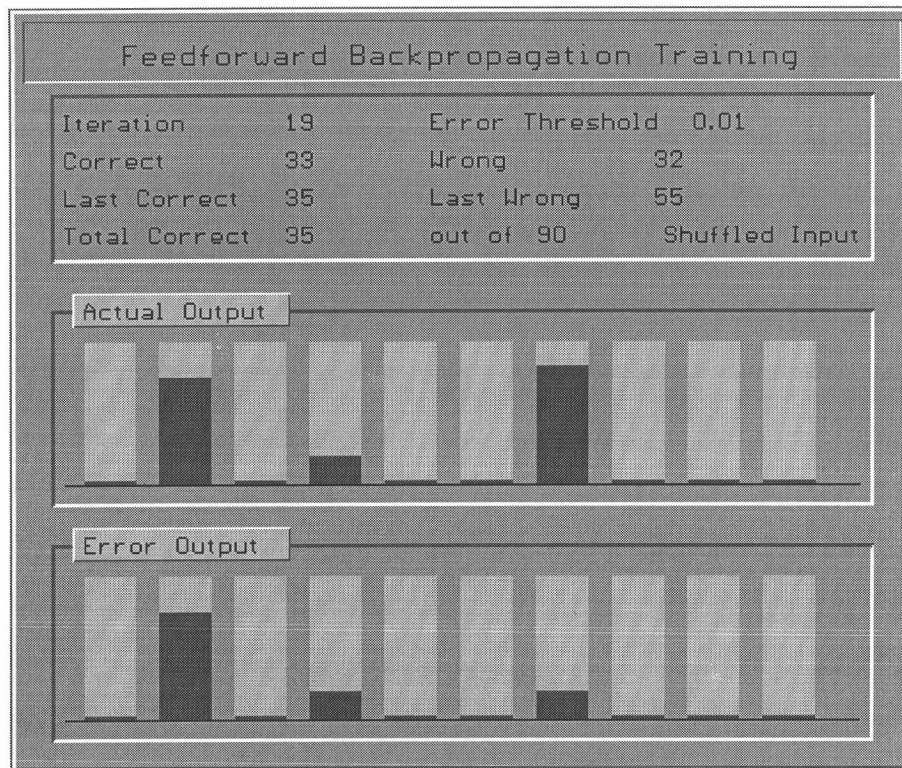
เนื่องจากกลุ่มตัวอักษร อ อ อ อ ไม่สามารถใช้รูปแบบตัวอักษรตัวหนาในการฝึกระบบได้ เนื่องจากภาพตัวอักษร อ และ อ ของรูปแบบตัวอักษรตัวหนาบางตัวบางขนาดจะมีแบบเปรียบ (primitive) ที่เหมือนกัน ซึ่งจะเป็นผลทำให้การฝึกหัดไม่สามารถบรรลุเข้าสู่ข้อผิดพลาดที่น้อยที่สุดได้ ดังนั้นเฉพาะชุดฝึกของภาพตัวอักษรชุดนี้จะใช้เฉพาะรูปแบบตัวอักษรปกติ และรูปแบบตัวอักษรแบบตัวเอนเท่านั้น ส่วนกลุ่มภาพตัวอักษรที่เหลือทั้งหมดจะใช้ภาพตัวอักษรรูปแบบตัวปกติ, ตัวเอน และตัวหนาในการฝึกหัดระบบ

สาเหตุที่ต้องแบ่งภาพตัวอักษรในการฝึกออกเป็นกลุ่มนั้น มีสาเหตุมาจากจำนวนตัวอักษรที่ใช้ในการฝึกมีจำนวนมาก (ภาพตัวอักษรของชุดฝึกทั้งหมด 690 ตัวอักษร) เพราะว่าจะต้องฝึกระบบให้รู้จำภาพตัวอักษรหลายขนาด (ภาพตัวอักษรขนาด 20, 24 และ 32 จุด) และหลายรูปแบบ (รูปแบบตัวปกติ, ตัวเอน และตัวหนา) ซึ่งทำให้หน่วยความจำของเครื่องคอมพิวเตอร์ไม่เพียงพอ และเวลาที่ใช้ในการฝึกจะใช้เวลานานมาก



รูปที่ 4.1 แสดงการสร้างข้อมูลชุดฝึก (training set) จากภาพตัวอักษร

จากนั้นจะนำภาพตัวอักษรแต่ละตัวของแต่ละกลุ่มมาสร้างแฟ้มแบบฝึก (training file) ของข้อมูลตัวอักษรแต่ละชุด โดยใช้ภาพข้อมูลตัวอักษร 9 ตัวอักษร (ภาพตัวอักษรขนาด 20, 24, และ 32 จุด รูปแบบตัวอักษรแบบตัวปกติ, ตัวเอน และ ตัวหนา) ในการฝึกตัวอักษรแต่ละตัว เมื่อสร้างแฟ้มแบบฝึก (training file) จากตัวอักษรเรียบร้อยแล้ว โปรแกรมจะนำแฟ้มแบบฝึกดังกล่าวมาสร้างแฟ้มน้ำหนัก (weight file) โดยที่โปรแกรมจะทำการคำนวณน้ำหนักการเชื่อมต่อ (weight connection) ของแต่ละโหนดแล้วจัดเก็บข้อมูลของการเชื่อมต่อลงในแฟ้มน้ำหนัก โดยมีเงื่อนไขของการฝึกระบบคือ ค่าความผิดพลาดที่ยอมรับได้ (error threshold) ในการทดลองกำหนดไว้ที่ 0.01, learning rate และ momentum มีค่า 1 และ 0.9 ตามลำดับ

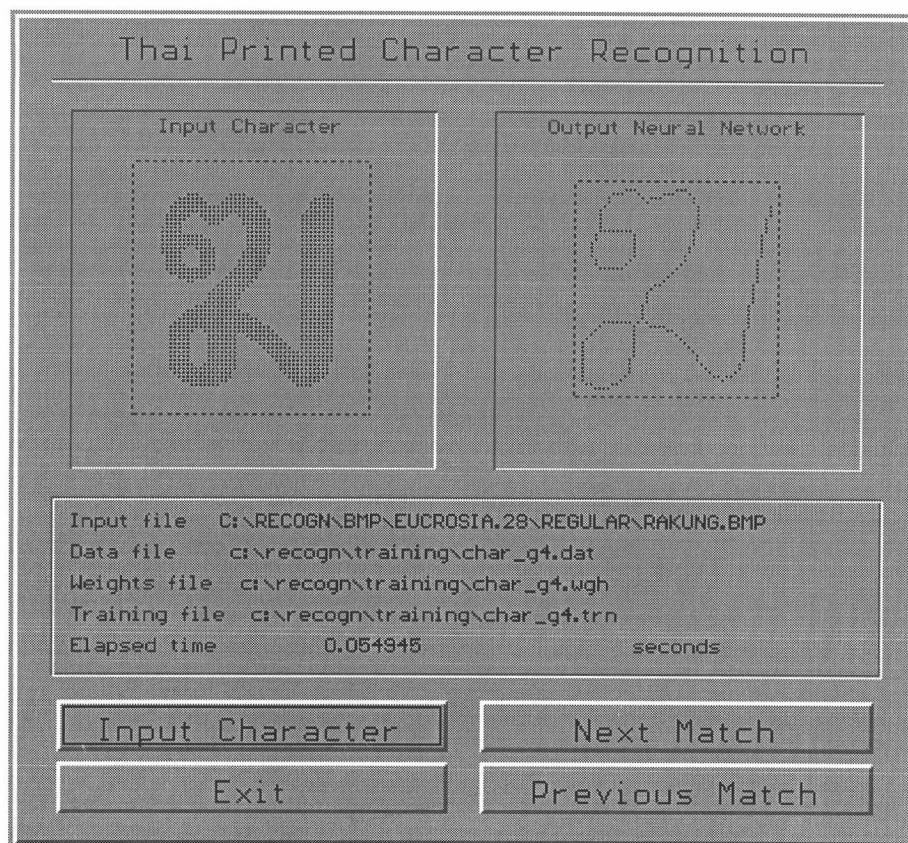


รูปที่ 4.2 แสดงการฝึกหัดเพื่อหาน้ำหนักการเชื่อมต่อ (weight connection)

การทดสอบการรู้จำ

เมื่อระบบนิรอลเน็ตเวิร์กเรียนรู้ภาพข้อมูลตัวอักษรจากชุดฝึกเรียบร้อยแล้ว ระบบนิรอลเน็ตเวิร์กก็จะสามารถรู้จำภาพตัวอักษรได้ ซึ่งในขั้นตอนแรกจะต้องอ่านค่าน้ำหนักการเชื่อมต่อ (weight connection) จากเพิ่มน้ำหนักที่ได้จากการฝึกหัดของระบบ โดยข้อมูลของภาพตัวอักษรแต่ละกลุ่มก็จะมีน้ำหนักการเชื่อมต่อแต่ละชุดเช่นกัน

ในการทดสอบนั้นจะทำการอ่านภาพข้อมูลของตัวอักษรแต่ละตัวอักษร หลังจากนั้นโปรแกรมจะแสดงตัวอักษรที่เป็นคำตอบซึ่งเก็บอยู่ในเพิ่มข้อมูล (data file) ของชุดฝึก ในการทดลองจะแบ่งข้อมูลภาพตัวอักษรเป็น 2 ชุด คือ ชุดของข้อมูลภาพตัวอักษรที่ระบบนิรอลเน็ตเวิร์กได้มีการเรียนรู้มาก่อน (ข้อมูลของชุดฝึก) และข้อมูลภาพตัวอักษรที่ระบบนิรอลเน็ตเวิร์กไม่เคยเรียนรู้มาก่อน (ข้อมูลของชุดทดสอบ)



รูปที่ 4.3 แสดงการทดสอบการรู้จำตัวอักษรภาษาไทย

การทดลองใช้เครื่องคอมพิวเตอร์ 486DX2-66 ในการทดสอบการรู้จำตัวอักษรภาษาไทย ซึ่งได้ผลการทดลองดังตารางที่ 4.1 ,ตารางที่ 4.2 และ ตารางที่ 4.3 มีอัตราการเรียนรู้ของภาพตัวอักษรที่ระบบเคยเรียนรู้มาก่อน 100% (690 ตัวอักษร) และมีอัตราการเรียนรู้จำตัวอักษรของภาพข้อมูลที่ระบบไม่เคยเรียนรู้มาก่อน 98.58% (702 ตัวอักษร) เวลาที่ใช้ในการประมวลผลเฉลี่ย 0.055 วินาทีต่อ 1 ภาพตัวอักษร ซึ่งเวลาที่ใช้ในการประมวลผลจะไม่ขึ้นอยู่กับรูปแบบความซับซ้อนของภาพตัวอักษร โดยมีภาพตัวอักษรที่รู้จำผิด 10 ภาพตัวอักษรดังแสดงในตารางที่ 4.4

ตารางที่ 4.3 แสดงผลการรู้จำตัวอักษรภาษาไทยที่รวมข้อมูลชุดฝึกและชุดทดสอบไว้ด้วยกัน

รูปแบบตัวอักษร	ขนาด (point)	จำนวนตัวอักษร	รู้จำถูกต้อง	รู้จำผิด	ความถูกต้อง (%)	เวลาเฉลี่ย (วินาที)
ตัวปกติ	20	78	78	0	100	0.055
	22	78	77	1	98.72	0.055
	24	78	78	0	100	0.055
	28	78	78	0	100	0.055
	32	78	78	0	100	0.055
	36	78	77	77	1	98.72
ผลรวม		468	466	2	99.57	0.055
ตัวเอน	20	78	78	0	100	0.055
	22	78	78	0	100	0.055
	24	78	78	0	100	0.055
	28	78	75	3	96.15	0.055
	32	78	78	0	100	0.055
	36	78	76	76	2	97.44
ผลรวม		468	463	5	98.93	0.055
ตัวหนา	20	74	74	0	100	0.055
	22	78	77	1	98.72	0.055
	24	74	74	0	100	0.055
	28	78	77	1	98.72	0.055
	32	74	74	0	100	0.055
	36	78	77	77	1	98.72
ผลรวม		456	453	3	99.34	0.055
ผลรวมทั้งหมด		1392	1382	10	99.28	0.055

ตารางที่ 4.3 แสดงผลการรู้จำตัวอักษรไทยที่รวมภาพข้อมูลตัวอักษรของชุดฝึก และชุดทดสอบเข้าไว้ด้วยกัน จะให้อัตราการรู้จำเฉลี่ยทั้งหมด 99.28% ตารางที่ 4.4 แสดงภาพตัวอักษรที่รู้จำผิด และตารางที่ 5.4 แสดงจำนวนโหนดของระดับข้อมูลเข้า (input node), จำนวนโหนดของระดับซ่อนตัว (hidden node), จำนวนโหนดของระดับแสดงผล (output node) และเวลาที่ใช้ในการฝึก

ตารางที่ 4.4 แสดงตัวอักษรที่รู้จำผิด

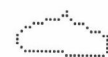
ตัวอักษรที่รู้จำผิด	ขนาด	ตัวปกติ	ตัวเอน	ตัวหนา	รู้จำผิดเป็น
ฉ	22	-	-	-	-
	28	-	1	-	ฉ
	36	-	-	-	-
ฉ	22	-	-	-	-
	28	-	-	-	-
	36	1	1	-	ฐ, ฒ
ภ	22	-	-	-	-
	28	-	1	-	ถ
	36	-	1	-	ถ
ว	22	-	-	-	-
	28	-	1	-	ร
	36	-	-	-	-
๙	22	1	-	1	๙
	28	-	-	1	๙
	36	-	-	1	๙
ผลรวมทั้งหมด		2	5	3	10/1392

จึงจำเป็นต้องแบ่งภาพข้อมูลตัวอักษรออกเป็นกลุ่มย่อย ๆ เพื่อที่จะลดขนาดของระดับแสดงผล (output layer) และเพิ่มจำนวนของโหนดในระดับซ่อนตัว (hidden layer)

การทำตัวอักษรให้บาง จะมีผลกระทบต่อภาพตัวอักษรที่มีลักษณะคล้ายคลึงกัน เช่น ข, ข, ค, ค, ช, ช, ญ, ญ, ท, ท ฯลฯ เนื่องจากเส้นหยักจะหายไปเมื่อผ่านกระบวนการทำภาพให้บาง ยกตัวอย่าง เช่น สระอ้อ เมื่อผ่านกระบวนการทำภาพให้บางแล้วเส้นตรงด้านซ้ายมือจะหายไปดังแสดง ในรูปที่ 4.4.ข ซึ่งทำให้มีรูปร่างลักษณะเหมือนกับสระอือ และเมื่อภาพดังกล่าวผ่านกระบวนการแปลงภาพข้อมูลเป็นเวกเตอร์ และเปลี่ยนเวกเตอร์เป็นแบบเปรียบเทียบจะได้ข้อมูลทางแบบเปรียบเทียบที่เหมือน กับสระ อือ ดังนั้นจะทำให้การรู้จำสระอือผิดไปและได้คำตอบเป็นสระอือแทน



รูปที่ 4.4.ก แสดงภาพก่อนการทำให้บาง



รูปที่ 4.4.ข แสดงภาพหลังการทำให้บาง

รูปที่ 4.4 แสดงภาพตัวอักษรสระอือที่มีปัญหาเมื่อผ่านการทำภาพตัวอักษรให้บาง

จากตารางที่ 4.1, 4.2, 4.3 แสดงผลการรู้จำตัวอักษรภาษาไทย และตารางที่ 4.4 แสดงภาพตัวอักษรที่ระบบรู้จำผิด จะพบว่าภาพตัวอักษรรูปแบบตัวเอนมีอัตราการรู้จำเฉลี่ยน้อยกว่าภาพตัวอักษรรูปแบบตัวปกติและตัวหนา โดยภาพตัวอักษรรูปแบบตัวเอนมีอัตราการรู้จำเฉลี่ย 98.93%, ภาพตัวอักษรรูปแบบตัวปกติมีอัตราการรู้จำเฉลี่ย 99.57% และภาพตัวอักษรรูปแบบตัวหนามีอัตราการรู้จำเฉลี่ย 99.34% สาเหตุที่ภาพตัวอักษรรูปแบบตัวเอนมีอัตราการรู้จำต่ำกว่าภาพตัวอักษรรูปแบบตัวปกติ และรูปแบบตัวหนา เนื่องจากตัวอักษรรูปแบบตัวปกติและตัวอักษรรูปแบบตัวหนา เมื่อผ่านกระบวนการการแปลงภาพข้อมูลให้เป็นแบบเปรียบเทียบจะได้ข้อมูลของแบบเปรียบเทียบที่คล้ายคลึงกันมาก กล่าวคือข้อมูลชุดฝึก (training set) ของภาพตัวอักษรรูปแบบตัวปกติ และตัวหนาจะมากกว่าตัวอักษรรูปแบบตัวเอนนั่นเอง