

ข้อมูลที่เลือกร้อยละ 10 นั้น จะทำการสุ่มมาทำข้อมูลในการสอนประมาณ 13,499 ชุด (Patterns) ทำการสกัดคุณลักษณะโดยใช้ขั้นตอนวิธีการวิเคราะห์องค์ประกอบหลัก (Principal Component Analysis) ซึ่งได้คุณลักษณะเด่นจากข้อมูลทั้งหมด 37 มิติ สำหรับการรู้จำด้วยโครงข่ายประสาทเทียมจำนวน 19 มิติ โดยเราจะพิจารณาจากค่าไอเกนของข้อมูล (จากที่กล่าวไว้ก่อนหน้านี้ว่าข้อมูล KDDcup99 เป็นข้อมูลทั้งหมด 41 มิติ แต่มิติที่เป็น Basic Features และมิติที่เป็นคำตอบจะไม่นำมาพิจารณา ดังนั้น จึงเหลือเพียง 37 มิติ)

4.2 การรู้จำประเภทของผู้บุกรุกเบื้องต้น

ในการทดลองเบื้องต้นสำหรับการรู้จำประเภทของผู้บุกรุกในงานวิจัยนี้ ผู้วิจัยเลือกวิธีการรู้จำแบบมีผู้สอน (Supervised learning) ที่ได้รับความนิยมในการใช้ทดสอบการรู้จำ คือ โครงข่ายประสาทเทียมแบบแพร่ย้อนกลับ และ Support Vector Machine โดยมีข้อมูลนำเข้าสำหรับการรู้จำ 2 ประเภท คือ ข้อมูลทั้งหมด 37 มิติ และ ข้อมูลที่ผ่านขั้นตอนการลดมิติข้อมูล (PCA) 19 มิติ ทำให้เราสามารถแบ่งการทดลองออกเป็น 4 การทดลอง ดังนี้

1. All+BPNN (ข้อมูลทั้งหมด 37 มิติ และรู้จำด้วยโครงข่ายประสาทเทียมแบบแพร่ย้อนกลับ)
 - Number of hiddenLayers = (attribs + classes) / 2
 - LearningRate=0.3
 - Momentum=0.2
 - TrainingTime=500
 - Training 100%
2. All+SVM (ข้อมูลทั้งหมด 37 มิติ และรู้จำด้วย Support Vector Machine)
 - The polynomial kernel
3. PCA+BPNN (ข้อมูลผ่าน PCA 19 มิติ และรู้จำด้วยโครงข่ายประสาทเทียมแบบแพร่ย้อนกลับ)
 - Number of hiddenLayers = (attribs + classes) / 2
 - LearningRate=0.3
 - Momentum=0.2
 - TrainingTime=500
 - Training 100%
4. PCA+SVM (ข้อมูลผ่าน PCA 19 มิติ และรู้จำด้วย Support Vector Machine)
 - The polynomial kernel

ตารางที่ 4-1 นำเสนอรายละเอียดของข้อมูลที่ใช้ในการทดลองนี้ โดยข้อมูลที่ใช้มี ชื่อ KDDcup99 ทำการสุ่มเลือกมาทั้งหมด 13,499 ชุดทดสอบ โดยข้อมูลมี 37 มิติ โดยรายละเอียดของข้อมูลแต่ละคลาสเป็นดังนี้

คลาสที่ 1 (ประเภทข้อมูล Normal) จำนวนข้อมูลที่สุ่มมาได้ 4,107 ชุดข้อมูล
 คลาสที่ 2 (ประเภทผู้บุกรุก DoS) จำนวนข้อมูลที่สุ่มมาได้ 4,107 ชุดข้อมูล
 คลาสที่ 3 (ประเภทผู้บุกรุก Probe) จำนวนข้อมูลที่สุ่มมาได้ 4,107 ชุดข้อมูล
 คลาสที่ 4 (ประเภทผู้บุกรุก R2L) จำนวนข้อมูลทั้งหมด 1,126 ชุดข้อมูล
 คลาสที่ 5 (ประเภทผู้บุกรุก U2L) จำนวนข้อมูลทั้งหมด 52 ชุดข้อมูล

ตารางที่ 4-1 รายละเอียดข้อมูลที่ใช้ในการทดลอง

ประเภทข้อมูล	จำนวนข้อมูล/มิติ (แอทริบิวต์)	จำนวนข้อมูล (Patterns) ในแต่ละคลาส
KDDcup99	13499/37	4107/4107/4107/1126/52

ตารางที่ 4-2 นำเสนอค่าสถิติที่ได้จากการทำการทดลองประกอบด้วย ค่าร้อยละของความถูกต้อง (Accuracy), ค่า F และ ค่า AUC พบว่า ในการทดลองเบื้องต้นในรายงานการวิจัยนี้ ชุดข้อมูลทั้งหมด 37 มิติ ให้ผลการทดลองที่สูงกว่าผลการทดลองที่ได้จากการลดมิติข้อมูลด้วยเทคนิค PCA

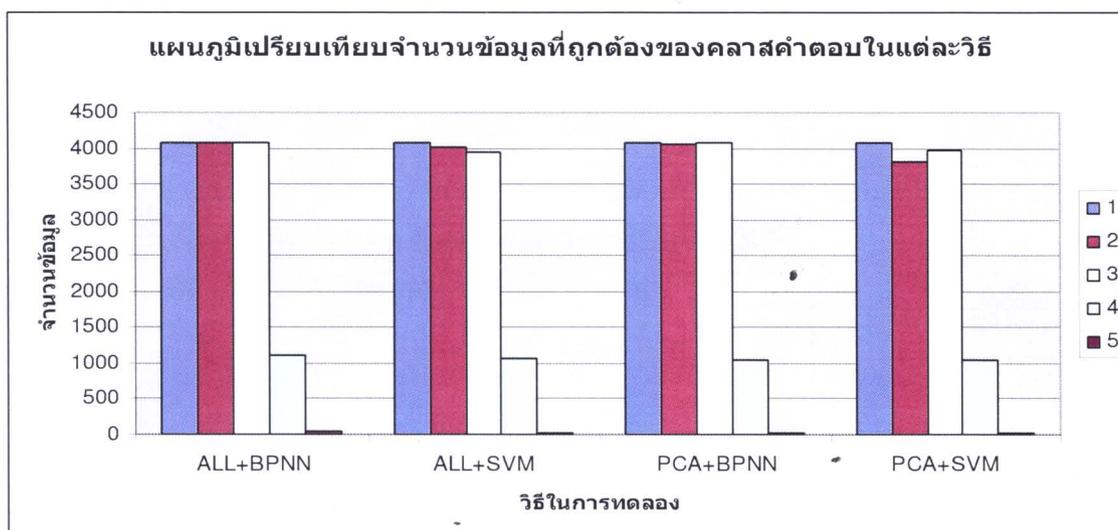
ตารางที่ 4-2 ค่า Accuracy จากการประมวลผล

Learning Method	ข้อมูลทั้งหมด			ข้อมูลที่ผ่านขั้นตอน PCA		
	Accuracy	F	AUC	Accuracy	F	AUC
BPNN	99.17	0.992	0.999	98.45	0.984	0.996
SVM	97.15	0.971	0.988	95.63	0.956	0.983

ตารางที่ 4-3 และ รูปที่ 4-2 ได้แสดงจำนวนชุดข้อมูลที่วิธีการเรียนรู้แต่ละวิธีทำการแบ่งประเภทได้อย่างถูกต้อง ซึ่งพบว่าในการทดลองเบื้องต้นนี้ โครงข่ายประสาทเทียมแบบแพร่ย้อนกลับ ให้ผลการทดลองที่ดีกว่า Support Vector Machine

ตารางที่ 4-3 จำนวนข้อมูลที่แบ่งประเภทได้ถูกต้องของคลาสคำตอบในแต่ละวิธี

Learning Method	คลาส				
	1 (4107)	2 (4107)	3 (4107)	4 (1126)	5 (52)
ALL+BPNN	4087	4073	4082	1112	34
ALL+SVM	4075	4002	3953	1057	28
PCA+BPNN	4072	4060	4080	1050	28
PCA+SVM	4075	3820	3957	1032	26



รูปที่ 4-2 แผนภูมิเปรียบเทียบจำนวนข้อมูลที่ถูกต้องของคลาสคำตอบในแต่ละวิธี