

การตรวจจับการบุกรุกโดยใช้ Self-Organizing Map หลายลำดับชั้น

INTRUSION DETECTION USING MULTI-LAYER
SELF-ORGANIZING MAP

สุรพล โรจนประดิษฐ์
SURAPOL ROCHANAPRATISHTHA

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรมหาบัณฑิต

สาขาวิชาวิศวกรรมคอมพิวเตอร์

บัณฑิตวิทยาลัย

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

พ.ศ. 2549

ISBN 974-15-2642-3

การตรวจจับการบุกรุกโดยใช้ Self-Organizing Map หลายลำดับชั้น

INTRUSION DETECTION USING MULTI-LAYER
SELF-ORGANIZING MAP



สุรพล โรจนประดิษฐ์

SURAPOL ROCHANAPRATISHTHA

เลขหมู่.....
เลขทะเบียน..... 63674
วัน,เดือน,ปี 30 ส.ค. 2549

.b.....
.i.....

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิศวกรรมศาสตรมหาบัณฑิต
สาขาวิชาวิศวกรรมคอมพิวเตอร์

บัณฑิตวิทยาลัย

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

พ.ศ.2549

ISBN 974-15-2642-3

**INTRUSION DETECTION USING MULTI-LAYER
SELF-ORGANIZING MAP**

SURAPOL ROCHANAPRATISHTHA

**A THESIS SUBMITTED IN PARTIAL FULFILLMENT
OF THE REQUIREMENT FOR THE DEGREE OF
MASTER OF ENGINEERING IN COMPUTER ENGINEERING
SCHOOL OF GRADUATE STUDIES
KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG**

2006

ISBN 974-15-2642-3

COPYRIGHT 2006

SCHOOL OF GRADUATE STUDIES

KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

บัณฑิตวิทยาลัย
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
ใบรับรองวิทยานิพนธ์

หัวข้อวิทยานิพนธ์ การตรวจจับการบุกรุกโดยใช้ SELF-ORGANIZING MAP หลายลำดับชั้น
INTRUSION DETECTION USING MULTI-LAYER SELF-ORGANIZING MAP


นักศึกษา นายสุรพล โรจนประดิษฐ

รหัสประจำตัว 45061217

ปริญญา วิศวกรรมศาสตรมหาบัณฑิต

สาขาวิชา วิศวกรรมคอมพิวเตอร์

อาจารย์ผู้ควบคุมวิทยานิพนธ์ รศ.ดร.เอื้อน ปิ่นเงิน

คณะกรรมการสอบวิทยานิพนธ์		ลายมือชื่อ
รศ.ดร.บุญธิรี	เกรือตราชู	
ดร.สมศักดิ์	วลัยรัชต์	
ดร.วัชระ	ฉัตรวิริยะ	
ผศ.เกียรติคุณ	เจียรนัยธนะกิจ	
รศ.ดร.เอื้อน	ปิ่นเงิน	

วัน / เดือน / ปี ที่สอบ 16 พฤษภาคม 2549 เวลา 11.30-13.30 น.

สถานที่สอบ ณ อาคาร 12 ชั้น ชั้น 3 (ห้อง E12-301)

บัณฑิตวิทยาลัยรับรองแล้ว

(ผศ.ดร.จารุวัตร เจริญสุข)
คณบดีบัณฑิตวิทยาลัย

วันที่.....14.....เดือน.....ก.พ.๒๕๔๙.....พ.ศ.....๒๕๔๙.....

หัวข้อวิทยานิพนธ์	การตรวจจับการบุกรุกโดยใช้ Self-Organizing Map หลายลำดับ ชั้น
นักศึกษา	นายสุรพล โรจนประดิษฐ์
รหัสนักศึกษา	45061217
ปริญญา	วิศวกรรมศาสตรมหาบัณฑิต
สาขาวิชา	วิศวกรรมคอมพิวเตอร์
พ.ศ.	2549
อาจารย์ผู้ควบคุมวิทยานิพนธ์	รศ.ดร. เอื้อน ปิ่นเงิน

บทคัดย่อ

การตรวจจับการบุกรุกเครือข่ายโดยอาศัยการแยกประเภทการบุกรุกที่ใช้ Self-Organizing Map (SOM) เป็นการจัดข้อมูลพฤติกรรมการบุกรุกที่ได้จากคุณลักษณะทางระบบเครือข่ายหลายมิติ โดยแปลงให้อยู่ในรูปของแผนภาพ SOM สองมิติ ซึ่งจะให้ข้อมูลการบุกรุกระบบเครือข่ายที่มีลักษณะคล้ายกันถูกจัดกลุ่มให้อยู่ในโหนดเดียวกัน แต่ในบางกรณีจะเกิดการซ้อนทับกันของข้อมูล กล่าวคือข้อมูลในโหนดเดียวกันอาจเกิดจากการบุกรุกเครือข่ายที่ต่างประเภทกัน ทำให้ไม่สามารถระบุประเภทของการบุกรุกระบบเครือข่ายได้อย่างชัดเจน ปัญหาหลักของแผนภาพที่มีข้อมูลประเภทการบุกรุกที่ทับซ้อนกันนั้น SOM แบบลำดับชั้นเดียวไม่สามารถแก้ปัญหาได้ ดังนั้นงานวิจัยนี้ได้นำเสนอการแยกประเภทการบุกรุกโดยใช้ SOM แบบหลายลำดับชั้น โดยแยกการทำงานออกเป็นสองขั้นตอน ขั้นตอนแรกเป็นการแยกประเภทการบุกรุกเครือข่ายจากข้อมูลเบื้องต้นโดยใช้ SOM แบบลำดับชั้นเดียว และขั้นตอนที่สองตรวจสอบ SOM ที่ได้ว่าโหนดใดมีอัตราการซ้อนทับของข้อมูลน้อยกว่าค่าที่กำหนดก็จะทำการแตกข้อมูลเฉพาะ โหนดนั้นออกเป็นอีกหนึ่งลำดับชั้น จากการทดลองพบว่าด้วยวิธีที่นำเสนอสามารถเพิ่มเปอร์เซ็นต์ detection rate และลดเปอร์เซ็นต์ false positive rate

Thesis Title	Intrusion Detection using Multi-layer Self-Organizing Map
Student	Mr. Surapol Rochanapratishtha
Student ID.	45061217
Degree	Master of Engineering
Programme	Computer engineering
Year	2006
Thesis Advisor	Assoc. Prof. Dr. Ouen Pinngern

ABSTRACT

The classification of network intrusion detection by using Self-Organizing Map (SOM) is an intrusion behavior management that mapped the multi-dimensional features into two dimensional SOM. Consequently, the similarities of intrusion data are classified in the same node. However, there are overlapping data in some cases such as data in a specific cluster have different types of attack. Therefore, it is difficult to identify the type of intrusion behavior. The problem of overlapping data can not be solved by single-layer SOM. This research presented the classification of network intrusion detection using multi-layer SOM in order to classify the types of intrusion. The process consists of two steps. First, the algorithm uses single-layer SOM to classify types of intrusion from the primary data. Second, the results from SOM are examined to determine the nodes that have overlapping rate less than the threshold. Then the data in one layer is distributed again. From the experiments, we found that the percentage of detection rate and false positive rate were improved.

กิตติกรรมประกาศ

วิทยานิพนธ์นี้สำเร็จลุล่วงด้วยดีเนื่องด้วยการอำนวยการที่ยิ่งใหญ่จากพระเจ้าพระเยซูคริสต์ และกำลังใจจาก คุณพ่อ คุณแม่ ข้าพเจ้าขอสำนึกในพระคุณนี้อย่างเป็นที่สุด

วิทยานิพนธ์นี้จะไม่สำเร็จลุล่วงหากปราศจากแรงผลักดัน และคำแนะนำที่มีประโยชน์ของ รศ.ดร. เอื้อน ปิ่นเงิน ผู้ควบคุมวิทยานิพนธ์ ข้าพเจ้าขอกราบขอบพระคุณเป็นอย่างสูง

ข้าพเจ้าขอกราบเท้า คุณครูและอาจารย์ทุกท่านตั้งแต่เล็กจนเติบโตใหญ่ ที่ได้มอบวิชาความรู้ให้แก่ข้าพเจ้า รวมทั้งคำสั่งสอนและอบรมให้ข้าพเจ้าเป็นคนดี ข้าพเจ้าขอกราบขอบพระคุณเป็นอย่างสูง

ข้าพเจ้าขอขอบคุณบัณฑิตวิทยาลัยที่ได้สนับสนุนเงินทุนการทำวิทยานิพนธ์และสำนักวิจัย การสื่อสารและเทคโนโลยีสารสนเทศ (ReCCIT) ที่ได้สนับสนุนเครื่องมือ ตลอดจนข้อมูล และหนังสือต่างๆ ที่ใช้ในการทำวิจัย

ข้าพเจ้าขอขอบคุณสำหรับกำลังใจ คำแนะนำ และประสบการณ์ที่ดีจากพี่ ๆ และเพื่อน ๆ นักศึกษา ป.โททุกท่าน และขอขอบคุณ นายพรเทพ โรจนวสุ นายไพฑูรย์ ศรีนิล นายกษานต์ ศรีกุล นาด นายธนวัฒน์ ภัทรวรรณ และนายณรงค์ชัย มุ่งแสงกลาง ที่ช่วยแก้ไขภาษาในการพิมพ์บทความตีพิมพ์ และวิทยานิพนธ์ฉบับนี้ ข้าพเจ้าขอขอบคุณ

สุดท้ายนี้ต้องขอขอบคุณ นางสาว พิมพ์ประไพ พุทธิवास ที่เป็นเสมือนเพื่อนคู่คิดและเป็นกำลังใจที่ดีตลอดมา

สำหรับคุณค่าและประโยชน์อันพึงมีจากวิทยานิพนธ์ฉบับนี้ ข้าพเจ้าขอมอบให้กับผู้มีพระคุณทุกท่าน หากวิทยานิพนธ์ฉบับนี้มีข้อผิดพลาดประการใดข้าพเจ้าขอน้อมรับไว้เพียงผู้เดียว

สุรพล โรจนประดิษฐ

สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	I
บทคัดย่อภาษาอังกฤษ.....	II
กิตติกรรมประกาศ.....	III
สารบัญ.....	IV
สารบัญตาราง.....	VII
สารบัญภาพ.....	VIII
บทที่ 1 บทนำ.....	1
1.1 ความเป็นมาและความสำคัญของปัญหา.....	1
1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา.....	1
1.3 แนวคิดที่ใช้ในงานวิจัย.....	2
1.4 การเปรียบเทียบระหว่างวิธีการที่นำเสนอกับวิธีการแบบพื้นฐาน.....	2
1.5 ขอบเขตการวิจัย.....	3
1.6 ขั้นตอนของการศึกษา.....	3
บทที่ 2 ระบบตรวจจับการบุกรุก.....	4
2.1 ความหมายของระบบตรวจจับการบุกรุก.....	4
2.2 ประเภทของระบบตรวจจับการบุกรุก.....	6
2.2.1 Anomaly Detection.....	6
2.2.2 Misuse Detection.....	7
2.3 การทำงานของระบบตรวจจับการบุกรุก.....	8
2.3.1 การเก็บข้อมูลในระบบ.....	8
2.3.2 การวิเคราะห์ข้อมูลระบบ.....	10
2.3.3 การตอบสนอง.....	13
2.3.4 การรายงานผลการทำงาน.....	14
2.4 ความสำคัญของระบบตรวจจับการบุกรุก.....	14
2.5 สรุป.....	15
บทที่ 3 นิเวศเน็ตเวิร์กแบบ ไม่มีผู้สอน.....	16
3.1 บทนำ.....	16

สารบัญ(ต่อ)

	หน้า
3.2 โมเดลและอัลกอริทึมของ SOM	16
3.3 คุณสมบัติของ SOM.....	22
3.3.1 แผนภาพเรียงตัว.....	22
3.3.2 การจัดแบ่งกลุ่มข้อมูล	23
3.4 งานวิจัยที่ประยุกต์ใช้ SOM.....	23
3.5 สรุป.....	27
บทที่ 4 การออกแบบมัลติเลเยอร์ SOM.....	28
4.1 โมเดลการทำงานของระบบ	28
4.2 ข้อมูล KDD Cup 1999	29
4.2.1 ลักษณะข้อมูลของ KDD Cup 1999	29
4.2.2 การเตรียมข้อมูลอินพุตเวกเตอร์.....	32
4.3 การสร้างแผนภาพ SOM ลำดับชั้นที่หนึ่ง.....	34
4.4 แผนภาพ SOM ลำดับชั้นต่อไป.....	35
4.5 อัลกอริทึมในการเพิ่มชั้นของแผนภาพ SOM	36
4.6 การวัดค่า Detection Rate และ False Positive	37
4.7 สรุป.....	38
บทที่ 5 ผลการทดลอง	40
5.1 การทดลองที่ 1 การตรวจจับการบุกรุกโดยใช้ SOM แบบลำดับชั้นเดียว	40
5.1.1 จุดประสงค์การทดลอง	40
5.1.2 ขั้นตอนการทดลอง	40
5.2 การทดลองที่ 2 การตรวจจับการบุกรุกโดยใช้ SOM แบบหลายลำดับชั้น.....	44
5.2.1 จุดประสงค์ของการทดลอง.....	44
5.2.2 ขั้นตอนการทดลอง	44
5.3 การทดลองที่ 3 การทดสอบกับข้อมูลด้วย Back propagation เปรียบเทียบกับ SOM แบบหลายลำดับชั้น.....	46
5.3.1 จุดประสงค์ของการทดลอง.....	46
5.3.2 ขั้นตอนการทดลอง.....	46

สารบัญ(ต่อ)

	หน้า
5.3 สรุป.....	48
บทที่ 6 สรุปผลการวิจัย และข้อเสนอแนะ	50
6.1 สรุปผลการวิจัย	50
6.1 ข้อเสนอแนะ	51
บรรณานุกรม.....	52
ภาคผนวก งานวิจัยที่ได้รับการตีพิมพ์	54
ประวัติผู้เขียน	62

สารบัญตาราง

ตารางที่	หน้า
3.1 แสดงอินพุตเวกเตอร์ในรูปแบบของ RGB.....	21
3.2 แสดงผลการทดลองของงานวิจัย.....	26
4.1 คุณลักษณะพื้นฐาน	29
4.2 Content features.....	30
4.3 Traffic feature.....	30
4.4 Host based feature	31
4.5 แสดงการกำหนดค่าตัวเลขแทน Protocol feature	32
4.6 แสดงการกำหนดค่าตัวเลขแทน Service feature.....	32
4.6 แสดงการกำหนดค่าตัวเลขแทน Service feature (ต่อ).....	33
4.7 แสดงการกำหนดค่าตัวเลขแทน Flag feature.....	33
5.1 พารามิเตอร์ในการสอนระบบเริ่มต้น	40
5.2 ข้อมูลที่ใช้ในการสอนและทดสอบระบบจำนวนทั้งสิ้น 492,843 เรคคอร์ด.....	41
5.3 ผลการทดสอบ	41
5.4 แสดงผลของการวัดค่า detection rate.....	43
5.5 แสดงผลของการวัดค่า false positive	43
5.6 แสดงผลของจำนวนโหนดในแผนภาพ SOM และจำนวนชั้นของแผนภาพ SOM.....	44
5.7 แสดงผลของการวัดค่า detection rate.....	45
5.8 แสดงผลของการวัดค่าfalse positive	45
5.9 แสดงรายละเอียดชุดข้อมูลมาตรฐาน	46
5.10 แสดงผลของจำนวนการจำแนกประเภทข้อมูล	47

สารบัญรูปภาพ

รูปที่	หน้า
2.1 ระบบการตรวจจับการบุกรุก.....	5
2.2 ขอบเขตที่เหลื่อมกันของ IDS กับระบบ ทำให้เกิด false positive และ false negative.....	6
2.3 Anomaly Detection Model	7
2.4 Misuse Detection Model.....	7
2.5 การทำงานของระบบตรวจจับการบุกรุก	8
3.1 แสดงโมเดลพื้นฐานของ SOM แบบสี่เหลี่ยม	17
3.2 แสดงโครงสร้างของ SOM.....	17
3.3 แสดงระยะทางแบบยูคลิดระหว่างเวกเตอร์ x และ m_j	18
3.4 แสดงกราฟของฟังก์ชัน Gaussian ($y=e^{-x^2}$).....	19
3.5 แสดงโครงสร้างของ SOM ขนาด 7×7	20
3.6 แสดงแผนภาพ SOM ณ จำนวนรอบที่แตกต่างกัน.....	21
3.7 แสดงตัวอย่างแผนภาพ SOM ขนาด 9×9 ณ จำนวนรอบที่แตกต่างกัน.....	22
3.8 แสดงแผนภาพ SOM จาก 140 เอกสาร โดยพื้นฐานข้อมูล LISA.....	23
3.9 แสดงแผนภาพ SOM ในงานวิจัย WEBSOM.....	24
3.10 แสดงการจัดกลุ่มเอกสาร โดยใช้แผนภาพ SOM ขนาด 10×15	24
3.11 แสดงชั้นวางหนังสือเสมือนใน LibViewer.....	25
3.12 แสดงโมเดลการทดลองของงานวิจัย.....	26
3.13 แสดงโมเดลการทดลองของงานวิจัย [6].....	27
4.1 แสดงขั้นตอนการทำงานของระบบ.....	28
4.2 แสดงแผนภาพ SOM ขนาด 900 โหนด หรือขนาด 30×30	34
4.3 แสดงตัวอย่างข้อมูลที่ทับซ้อนในแผนภาพ SOM.....	35
4.4 แสดงการแบ่งชั้นของแผนภาพ SOM หลายลำดับชั้น	36
5.1 แสดงแผนภาพ SOM ขนาด 900 โหนด หรือขนาด 30×30	42
5.2 แสดงตัวอย่างผลการทดลองกับชุดข้อมูล Wine.....	47
5.2 แสดงเปอร์เซ็นต์การตรวจจับ Detection rate	49
5.3 แสดงเปอร์เซ็นต์ False positive.....	49

บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

การตรวจจับการบุกรุกเป็นส่วนสำคัญส่วนหนึ่งในการรักษาความปลอดภัยบนระบบเครือข่ายคอมพิวเตอร์ แต่เนื่องจากในปัจจุบันการพัฒนาทางด้านระบบเครือข่ายคอมพิวเตอร์ได้มีการพัฒนาไปอย่างรวดเร็วและมีการนำเสนอต่อสาธารณะมากขึ้น จึงเกิดจํานวนรูปแบบการบุกรุกมากขึ้นตามไปด้วย ดังนั้นความสนใจในการพัฒนาระบบการตรวจจับผู้บุกรุกจึงมุ่งไปที่การใช้เทคนิคใดมาหนึ่งมาประยุกต์ใช้ในการตรวจจับผู้บุกรุกระบบเครือข่าย โดยเทคนิคหนึ่งที่ถูกนำมาใช้คือ เซลฟออร์แกนไนซิงแมป (Self-Organizing Map) หรือ SOM ซึ่งเป็นนิวรอลเน็ตเวิร์กแบบไม่มีผู้สอน (unsupervised neural network) ถูกนำเสนอโมเดลขึ้นในปี ค.ศ. 1982 โดยศาสตราจารย์โคโฮเฮน มีคุณสมบัติที่สำคัญคือ สามารถแสดงผลข้อมูลที่มีมิติสูงให้อยู่ในรูปแบบของแผนภาพสองมิติ หลังจากผ่านกระบวนการเรียนรู้แล้ว แผนภาพที่ได้จะอยู่ในลักษณะของแผนภาพจัดเรียงตัว (ordered map) กล่าวคือข้อมูลที่ใกล้เคียงกันจะถูกจัดลงในแผนภาพบริเวณใกล้เคียงกัน ช่วยในการวิเคราะห์คุณลักษณะของข้อมูลได้ เช่น การกระจาย ความหนาแน่น และความสัมพันธ์ของข้อมูล แต่การประยุกต์ใช้ SOM ดังกล่าวยังมีปัญหา คือ ยังมีความสามารถในการกระจายข้อมูลพฤติกรรมการบุกรุกไม่ดีพอ เนื่องจากมีข้อมูลพฤติกรรมการบุกรุกบางชนิดนั้นมีการทับซ้อนกันกับพฤติกรรมปกติ ทำให้ระบบการตรวจจับผู้บุกรุกไม่สามารถแยกประเภทการบุกรุกได้อย่างมีประสิทธิภาพ

แนวทางในการแก้ไขปัญหาดังกล่าวคือการพัฒนาประสิทธิภาพให้ระบบตรวจจับผู้บุกรุกที่ประยุกต์ใช้ SOM ให้มีประสิทธิภาพในการตรวจจับพฤติกรรม ได้ถูกต้องมากขึ้น โดยงานวิจัยนี้ นำเสนอวิธีการตรวจจับการบุกรุกโดยอาศัยการแยกประเภทข้อมูลการบุกรุกระบบเครือข่ายคอมพิวเตอร์ โดยแสดงให้อยู่ในรูป SOM แบบหลายลำดับชั้น โดยมีการคำนวณค่าอัตราส่วนค่าหนึ่งที่ได้จากการพิจารณาโหนดที่มีข้อมูลที่ทับซ้อนกันเพื่อใช้เป็นข้อกำหนดในการเพิ่มลำดับชั้น และทำการทดลองบนพื้นฐานการบุกรุกระบบเครือข่ายแบบ Denial of Service กับ Probing และพฤติกรรมปกติ ด้วยวิธีการนี้เราสามารถที่จะได้ประสิทธิภาพการตรวจจับผู้บุกรุกที่ใกล้เคียงกับโมเดลเดิม นอกจากนั้นค่าความผิดพลาดยังได้ค่าที่ดีกว่าโมเดลเดิม

1.2 ความมุ่งหมายและวัตถุประสงค์ของการศึกษา

1. เพื่อศึกษาการ โมเดล SOM อัลกอริทึม การเรียนรู้ของ SOM และขีดจำกัดของโมเดล SOM

2. เพื่อศึกษาแนวทางในการพัฒนา SOM และนำเสนอโมเดลใหม่ที่ผู้วิจัยได้พัฒนาขึ้น
3. เพื่อศึกษาการจัดกลุ่มข้อมูลการบุกรุกระบบเครือข่ายและการประยุกต์ใช้ SOM ในการจัดกลุ่มข้อมูลการบุกรุกระบบเครือข่าย
4. เพื่อเปรียบเทียบประสิทธิภาพโมเดล SOM แบบเดิมและโมเดลใหม่ที่ผู้วิจัยนำเสนอ

1.3 แนวคิดที่ใช้ในงานวิจัย

จุดอ่อนของซิงเกิลเลเยอร์เซลฟี่ออ์แกนไนซซิงแม็พ (Single-layer Self-Organizing Map) ในการจัดกลุ่มข้อมูลพฤติกรรมกรการบุกรุกระบบเครือข่าย คือการซ้อนทับกันของข้อมูลพฤติกรรมกรการบุกรุก ดังนั้นการตรวจจับพฤติกรรมกรการบุกรุกระบบเครือข่าย จึงไม่มีประสิทธิภาพ และไม่สามารถระบุได้ว่าข้อมูลนั้นเป็นการบุกรุกประเภทใด

การแก้ปัญหาข้างต้นนี้ เราใช้มัลติเลเยอร์เซลฟี่ออ์แกนไนซซิงแม็พ (Multi-layer Self-Organizing Map) ในการปรับปรุงประสิทธิภาพของการตรวจจับพฤติกรรมกรการบุกรุกระบบเครือข่าย โดยมีการนำเสนอการคำนวณค่าอัตราส่วนค่าหนึ่งเพื่อใช้ในการพิจารณาโหนดที่มีข้อมูลที่ซ้อนทับกัน และใช้เป็นข้อกำหนดในการพิจารณาการเพิ่มลำดับชั้นของเซลฟี่ออ์แกนไนซซิงแม็พ (Self-Organizing Map)

เนื่องจากปัญหาของการซ้อนทับกันของข้อมูลพฤติกรรมกรการบุกรุก จึงมีแนวคิดวิธีการเพิ่มลำดับชั้นของเซลฟี่ออ์แกนไนซซิงแม็พ เพื่อแก้ไขปัญหาดังกล่าว โดยลักษณะเด่นของวิธีการที่นำเสนอคือ จะทำการเพิ่มลำดับชั้นเฉพาะในโหนดที่มีข้อมูลที่ทับซ้อนกันเท่านั้น สามารถทำให้ปัญหาเรื่องการทับซ้อนกันของข้อมูลพฤติกรรมกรการบุกรุกได้รับการแก้ไข ในวิทยานิพนธ์นี้จะนำเสนออัลกอริทึมของการเพิ่มลำดับชั้น โดยใช้อัตราส่วนของกลุ่ม (Class ratio) ในการพิจารณาการเพิ่มลำดับชั้น และแสดงผลการตรวจจับพฤติกรรมกรการบุกรุกของมัลติเลเยอร์เซลฟี่ออ์แกนไนซซิงแม็พ (Multi-layer Self-Organizing Map) เพื่อเปรียบเทียบวิธีการตรวจจับการพฤติกรรมกรการบุกรุกของซิงเกิลเลเยอร์เซลฟี่ออ์แกนไนซซิงแม็พ (Single-layer Self-Organizing Map)

1.4 การเปรียบเทียบระหว่างวิธีการที่นำเสนอกับวิธีการแบบพื้นฐาน

วิธีการเปรียบเทียบประสิทธิภาพของการตรวจจับการบุกรุกระบบเครือข่าย โดยใช้สมการวัดค่าความถูกต้องในการตรวจจับ (detection rate) และสมการวัดค่าความผิดพลาดในการตรวจจับ (false positive rate) เมื่อเทียบกับหลักการในแบบพื้นฐานแล้ว ในส่วนของค่าความถูกต้องในการตรวจจับ จะให้ค่าความถูกต้องที่สูงกว่า และจะให้ค่าความผิดพลาดในการตรวจจับที่ต่ำกว่า

1.5 ขอบเขตการวิจัย

1. ศึกษาเปรียบเทียบประสิทธิภาพการจัดพฤติกรรมกรการบุกรุกของโมเดล Single-layer SOM กับ Multi-layer SOM
2. ข้อมูลพฤติกรรมกรการบุกรุกที่ใช้ในการทดสอบเป็นข้อมูลที่ได้จากการจำลองพฤติกรรมกรการบุกรุกบนระบบเครือข่าย โดยเลือกเฉพาะพฤติกรรมกรการบุกรุกประเภท Denial of Service กับ Probing และพฤติกรรมกรปกติ จากข้อมูลจำลองพฤติกรรมกรการบุกรุกชุดมาตรฐาน (KDD Cup 1999)

1.6 ขั้นตอนของการศึกษา

วิทยานิพนธ์ฉบับนี้ได้แบ่งเนื้อหาออกเป็น 6 บทด้วยกันคือ

บทที่ 1 กล่าวถึงความเป็นมาของงานวิจัย ความมุ่งหมายและวัตถุประสงค์ ขอบเขตของการวิจัย และขั้นตอนการศึกษา

บทที่ 2 กล่าวถึงระบบตรวจจับผู้บุกรุก ความหมายของระบบตรวจจับผู้บุกรุก ประเภทของระบบ การทำงานของระบบ และความสำคัญของระบบตรวจจับผู้บุกรุก

บทที่ 3 กล่าวถึงนิรอลเน็ตเวิร์คแบบไม่มีผู้สอน โดยจะเน้นที่โมเดลของโคโฮเนนที่มีชื่อว่า Self-Organizing Map (SOM) ซึ่งเป็น โมเดลที่งานวิจัยนี้นำมาใช้ในระบบตรวจจับผู้บุกรุก

บทที่ 4 กล่าวถึงหลักการและวิธีการดำเนินงานวิจัย โดยจะเน้นถึงการเตรียมระบบเพื่อทำการทดสอบ

บทที่ 5 เป็นการทดลอง ผลการทดลอง

บทที่ 6 เป็นบทสรุปผลการวิจัยและข้อเสนอแนะ

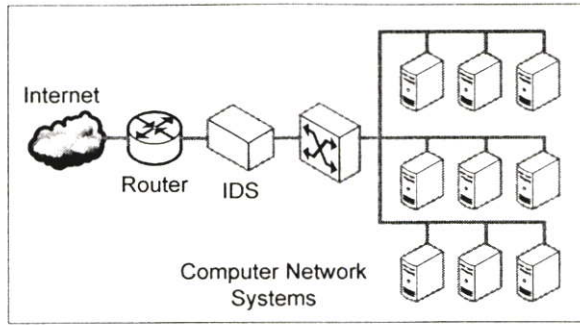
บทที่ 2

ระบบตรวจจับการบุกรุก

ปัจจุบันมีการใช้งานระบบคอมพิวเตอร์อย่างแพร่หลาย หลากๆ บริษัท หรือหน่วยงานต่างๆ ได้นำระบบคอมพิวเตอร์มาใช้เพื่อพัฒนาศักยภาพการทำงานของตนให้มากขึ้น การทำงานของระบบคอมพิวเตอร์จะทำงานอย่างต่อเนื่อง และมีการออนไลน์ให้คนอื่นๆ เข้ามาใช้งานระบบบางส่วนด้วย ซึ่งปัญหาที่เกิดขึ้นตามมาก็คือปัญหาการบุกรุกเข้ามาสร้างความเสียหายให้กับระบบ และอาจเข้ามาเพื่อขโมยข้อมูลที่สำคัญไป ปัญหาดังกล่าวจะเป็นปัญหากับผู้ดูแลระบบอย่างมาก เนื่องจากผู้ดูแลระบบต้องคอยป้องกันและแก้ไขปัญหาดังกล่าว อยู่เสมอๆ การทำงานของผู้ดูแลระบบนี้จะหนักมากหรือน้อยก็ขึ้นอยู่กับขนาดของระบบว่ามีขนาดใหญ่ และซับซ้อนมากน้อยเพียงใด ยิ่งระบบที่มีความซับซ้อนสูง มีความยุ่งยากในการดูแลมาก มีการออนไลน์ใช้งานอยู่ตลอดเวลาผู้ดูแลระบบก็จะต้องอยู่ดูแลระบบตลอดเวลาซึ่งในทางปฏิบัติแล้วเป็นไปได้แน่นอน อีกทั้งปัญหาบางปัญหาผู้ดูแลระบบไม่สามารถตรวจสอบด้วยตาของตัวเองได้เลย ปัญหาการดูแลระบบจึงยุ่งยากขึ้นทุกวัน ทางออกที่จะช่วยแก้ปัญหาเหล่านี้ ก็คือการมีผู้ช่วยมาช่วยดูแลระบบตลอดเวลาสามารถพบเห็นในสิ่งที่ผู้ดูแลระบบมองไม่เห็น คอยแจ้งเตือนให้กับผู้ดูแลระบบมีเหตุการณ์ผิดปกติเกิดขึ้น และในบางครั้งก็อาจสามารถแก้ไขปัญหาเบื้องต้น ได้ด้วย ผู้ช่วยที่สามารถช่วยแบ่งเบาภาระของผู้ดูแลระบบได้อย่างมากก็คือ ระบบตรวจจับการบุกรุก หรือ Intrusion Detection System

2.1 ความหมายของระบบตรวจจับการบุกรุก

ระบบตรวจจับการบุกรุก คือ ระบบตรวจจับสัญญาณของความผิดปกติต่างๆ ที่เกิดขึ้นในระบบที่อยู่ในขอบเขตที่ระบบนี้มีหน้าที่ตรวจสอบ ในที่นี้จะหมายถึง โปรแกรมที่ใช้สำหรับตรวจจับความผิดปกติในระบบเครือข่ายคอมพิวเตอร์เท่านั้น โดยตัวโปรแกรมจะมีความสามารถในการตรวจจับสัญญาณของความผิดปกติที่เกิดขึ้นในระบบ ไม่ว่าจะเป็นภายในระบบคอมพิวเตอร์ ระบบปฏิบัติการ โปรแกรมที่รันอยู่ในเครื่อง การทำงานกับฐานข้อมูล หรือแม้แต่ข้อมูลที่วิ่งผ่านไปมาในเครือข่ายด้วย

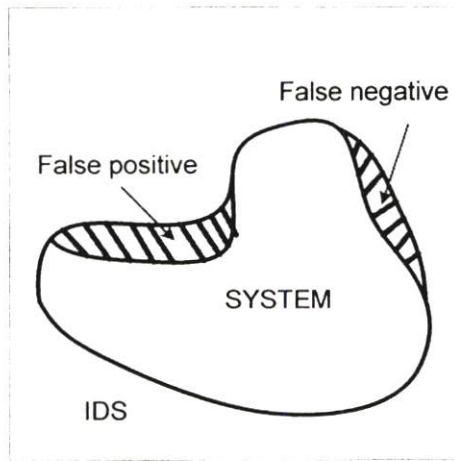


รูปที่ 2.1 ระบบการตรวจจับการบุกรุก

จากรูปที่ 2.1 ถ้าเรามองระบบเป็นเซตของการทำงานเซตหนึ่ง ระบบตรวจจับการบุกรุกที่แท้จริง (Ideal IDS) ต้องทราบขอบเขตของระบบว่าระบบทำงานอะไรบ้าง การทำงานใดปกติและการทำงานใดผิดปกติ โดยระบบตรวจจับการบุกรุกที่มีประสิทธิภาพ จะทราบขอบเขตของระบบ โดยไม่มีการเหลื่อมล้ำเข้าไปในระบบ หรือเหลื่อมล้ำออกนอกระบบ อย่างเด็ดขาด แต่ในระบบที่ใช้งานในโลกของความเป็นจริง กรอบของระบบที่ IDS รับรู้อาจมีความเหลื่อมล้ำกับระบบที่ IDS ต้องตรวจสอบ ทำให้เกิดความผิดพลาดในการตรวจสอบได้ ซึ่งสามารถแบ่งความผิดพลาด ในการตรวจสอบได้เป็นสองลักษณะคือ false positive และ false negative ดังรูปที่ 2.2 ซึ่งความผิดพลาดทั้งสองแบบมีรายละเอียดคือ

1. False Positive คือ ความผิดพลาดอันเนื่องมาจากการเกิดเหตุการณ์ซึ่งเป็นเหตุการณ์ปกติในระบบ แต่ IDS คิดว่าเกิดเหตุการณ์ผิดปกติเกิดขึ้น ผลลัพธ์คือ IDS จึงแจ้งเตือนผู้ดูแลระบบว่าเกิดเหตุการณ์ผิดปกติ
2. False Negative คือ ความผิดพลาดอันเนื่องมาจากการเกิดเหตุการณ์ซึ่งเป็นเหตุการณ์ผิดปกติในระบบ แต่ IDS คิดว่าเป็นเหตุการณ์ปกติ จึงไม่ได้แจ้งเตือนผู้ดูแลระบบว่าเกิดเหตุการณ์ผิดปกติ

โดยการออกแบบระบบ IDS นั้นจะพยายามออกแบบ และใช้วิธีต่างๆ มากมายที่ทำให้ False positive และ False negative มีน้อยที่สุดซึ่งก็มีงานวิจัยหลายๆ ชิ้นที่ทำเพื่อลดพื้นที่ตรงส่วนนี้ [1], [2], [3], [4], [5], [6]



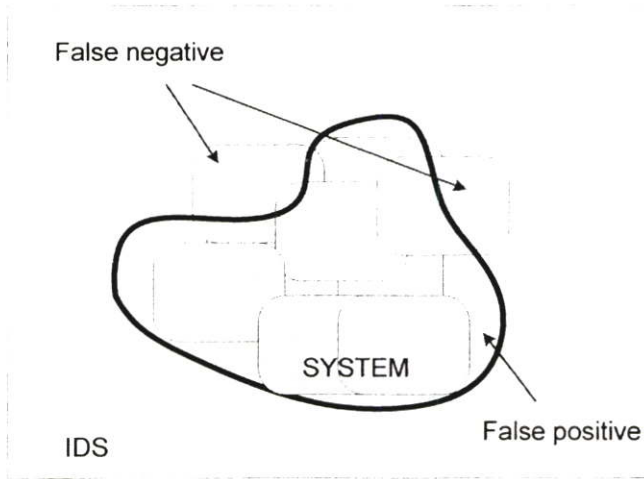
รูปที่ 2.2 ขอบเขตที่เหลื่อมกันของ IDS กับระบบ ทำให้เกิด false positive และ false negative

2.2 ประเภทของระบบตรวจจับการบุกรุก

ระบบตรวจจับการบุกรุกแบ่งออกเป็นสองรูปแบบ [7] คือ ระบบที่ตรวจหาการทำงานที่ผิดไปจากการทำงานปกติของระบบ เรียกว่า Anomaly Detection ซึ่งเป็นเหมือนกับการตรวจจับคนที่ไม่มีสิทธิทำงานอยู่ในระบบ อีกรูปแบบหนึ่ง คือ ระบบที่ตรวจหาการทำงานที่ไม่ควรเกิดขึ้นในระบบ เรียกว่า Misuse Detection ในที่นี้เปรียบเสมือนเป็นบุคคลที่มีสิทธิในระบบ สามารถเข้าออกในระบบได้ แต่เป็นผู้ที่ทำในสิ่งที่ระบบไม่อนุญาตให้ทำ หรือทำการใดๆ ที่อยู่นอกเหนือสิทธิของคนในระบบ

2.2.1 Anomaly Detection

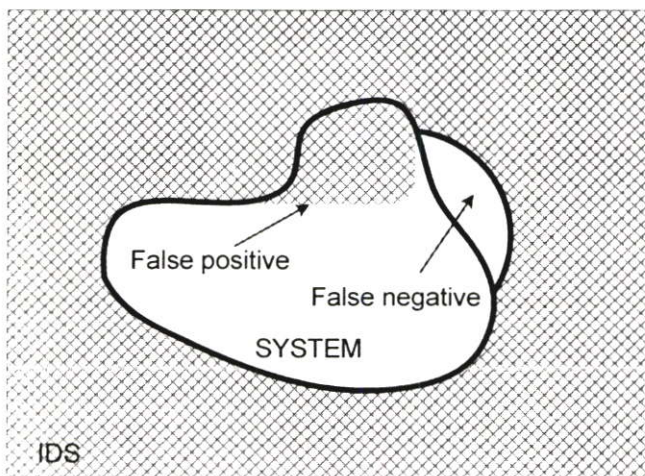
แนวความคิดของการทำ Anomaly Detection คือ การหาเซตของการทำงานที่เป็นปกติย่อยๆ ขึ้นมาแล้วนำมารวมกันเพื่อให้ระบบ IDS ทราบข้อมูลของเซตการทำงานที่เป็นปกติทั้งหมดในระบบ หลังจากนั้นเมื่อให้ระบบ IDS ทำงาน ถ้าเกิดกรณีที่ IDS ตรวจจับการทำงานที่ไม่ได้อยู่ในเซตของการทำงานที่เป็นปกติ ระบบ IDS จะแจ้งเตือนต่อผู้ดูแลระบบทันที สำหรับการสร้างขอบเขตของระบบนั้น อาจสร้างได้โดยการหาข้อมูลการทำงานที่เป็นปกติในระบบขึ้นมา โดยเอาข้อมูลการทำงานของผู้ใช้งานแต่ละคน เวลาที่มีการใช้งาน ทรัพยากรที่ผู้ใช้งานคนนั้นๆ มักจะใช้บ่อยๆ หรือแม้กระทั่งข้อมูลในระบบ หรือในเครือข่ายก็สามารถนำมาสร้างเป็น เซตของระบบได้เช่นกันดังตัวอย่างรูปที่ 2.3 ในการหาเซตของการทำงานที่เป็นปกติทั้งหมด อาจเกิดการผิดพลาดขึ้นมาทำให้เกิดลักษณะของ false positive และ false negative ขึ้นมาได้เช่นกัน



รูปที่ 2.3 Anomaly Detection Model

2.2.2 Misuse Detection

เป็นแนวความคิดที่ตรงข้ามกับ Anomaly Detection คือ รูปแบบนี้จะใช้ข้อมูลของการทำงานที่ผิดปกติต่างๆ ที่เคยเกิดขึ้นมาแล้ว สร้างเป็นฐานข้อมูลของการทำงานที่ผิดปกติให้ระบบ IDS จดจำไว้ และในการทำงานของ IDS ที่มีการทำงานแบบ Misuse จะนำข้อมูลที่อยู่ในระบบมาค้นหาในฐานข้อมูลว่ามีอยู่หรือไม่ ถ้าระบบ IDS มีข้อมูลของการทำงานรูปแบบนั้นๆ อยู่ ก็แสดงว่าเกิดความผิดปกติขึ้นแล้ว ดังรูปที่ 2.4 แต่ในการรวบรวมนี้อาจรวมเอาการทำงานที่เป็นปกติเข้าไปด้วย ทำให้เกิด false positive หรือในบางกรณีที่ไม่ได้เก็บข้อมูลความผิดปกติไว้ก็ทำให้เกิดกรณีของ false negative ได้เช่นกัน ซึ่งในการทำงานของ Misuse Detection นี้ จะมีข้อเสียคือจะไม่สามารถตรวจจับการบุกรุกชนิดใหม่ๆ ได้ เนื่องจากต้องมีข้อมูลของการบุกรุกอยู่ก่อน จึงจะตรวจจับได้

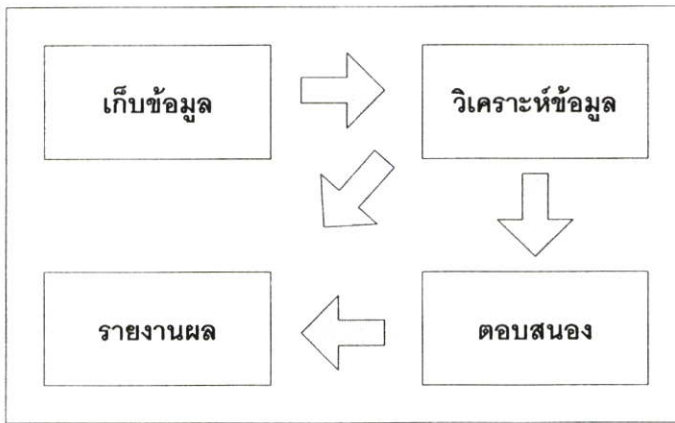


รูปที่ 2.4 Misuse Detection Model

2.3 การทำงานของระบบตรวจจับการบุกรุก

ระบบตรวจจับการบุกรุกแต่ละแบบมีหน้าที่การทำงานที่แตกต่างกันออกไป บางตัวจะตรวจจับความผิดปกติในระบบเครือข่าย บางตัวจะตรวจจับความผิดปกติในระบบฐานข้อมูล แต่โดยการทำงานทั้งหมดแล้วเราสามารถแบ่งการทำงานของระบบตรวจจับการบุกรุกได้เป็น 3 ขั้นตอน คือ การเก็บข้อมูลระบบ การวิเคราะห์ข้อมูลที่ได้ และการรายงานผลการทำงานให้ผู้ดูแลระบบหรือผู้ที่เกี่ยวข้องทราบ

ในการทำงานหลักๆ ของระบบตรวจจับการบุกรุก อาจมีขั้นตอนเสริมอยู่ขั้นตอนหนึ่ง คือ การตอบสนองต่อการบุกรุกนั้นๆ การทำงานในขั้นตอนนี้จะใช้ในกรณีที่การบุกรุกเป็นรูปแบบการบุกรุกที่ระบบตรวจจับการบุกรุกสามารถแก้ไขได้ด้วยตัวเองได้ ซึ่งในบางระบบอาจไม่มีการทำงานในส่วนนี้ ระบบที่ไม่มีการตอบสนองต่อการบุกรุก ส่วนใหญ่จะเป็นระบบที่ไม่ได้ทำงานแบบตอบสนองทันที (real time) คือ จะเก็บข้อมูลของระบบไว้ก่อน แล้วจึงวิเคราะห์ข้อมูลภายหลัง เมื่อทำการวิเคราะห์ข้อมูลแล้วพบว่ามี การบุกรุกเข้าสู่ระบบก็จะทำการแจ้งเตือนในขั้นตอนการทำการรายงานผลการทำงาน ระบบตรวจจับการบุกรุกที่ไม่มีการตอบสนองต่อการบุกรุกก็มักใช้ในงานที่ไม่มีความสำคัญมากนัก แต่ต้องการความถูกต้องสูง โดยลำดับการทำงานของระบบตรวจจับการบุกรุกสามารถมองเป็นขั้นตอนต่างๆ ได้ดังรูปที่ 2.5



รูปที่ 2.5 การทำงานของระบบตรวจจับการบุกรุก

2.3.1 การเก็บข้อมูลในระบบ

จากที่ได้กล่าวมาแล้วว่าระบบตรวจจับการบุกรุกจะมีการทำงานที่แตกต่างกันไป หน้าที่ในการเก็บข้อมูลของระบบที่ต้องการตรวจสอบก็แตกต่างกันไปตามหน้าที่ของระบบตรวจจับผู้บุกรุกด้วย โดยเราสามารถแบ่งการเก็บข้อมูลของระบบที่ต้องการตรวจสอบออกเป็นกลุ่มต่างๆ ได้ 4 กลุ่มด้วยกันคือ มีการเก็บข้อมูลในชั้นแอปพลิเคชัน (Application-based Approach) เพื่อนำมาตรวจสอบการทำงานของแอปพลิเคชันต่างๆ ว่าผิดปกติหรือไม่ การเก็บข้อมูลของการทำงานของเครื่อง (Host-based Approach) เพื่อนำมาตรวจสอบการทำงานของระบบของเครื่องที่

ใช้งานอยู่ การเก็บข้อมูลการเปลี่ยนแปลงข้อมูลในระบบ (Target-based Approach) เพื่อนำมาตรวจสอบว่าข้อมูลมีการเปลี่ยนแปลงอย่างไร และการเก็บข้อมูลเครือข่าย (Network-based Approach) เพื่อนำมาตรวจสอบว่ามีการบุกรุกทางระบบเครือข่ายหรือไม่ อย่างไร

2.3.1.1 การเก็บข้อมูลในชั้นแอปพลิเคชัน

การเก็บข้อมูลในชั้นแอปพลิเคชันนั้น เป็นการเก็บข้อมูลที่โปรแกรมต่างๆ สร้างขึ้นมาเพื่อรายงานผลการทำงานของโปรแกรมนั้นๆ เช่น log file หรือ error message ต่างๆ ของเว็บเซิร์ฟเวอร์ ไฟร์วอลล์ หรือ โปรแกรมบริหารฐานข้อมูล รวมถึงข้อมูลของการทำงานตอบสนองกันระหว่างผู้ใช้งาน โปรแกรม และข้อมูลที่เกี่ยวข้อง ในการเก็บข้อมูลในลักษณะนี้ นอกจากเป็นการทำงานด้านการรักษาความปลอดภัยในระบบแล้ว ยังช่วยในการวิเคราะห์ระบบและปรับปรุงระบบเนื่องจากผลจากการวิเคราะห์ข้อมูลที่ได้ทำให้ทราบว่าการทำงานของโปรแกรมไหนในระบบมีมากน้อยอย่างไร และควรให้ความสำคัญกับการทำงานตรงส่วนไหน แต่การทำงานในส่วนนี้ก็ยังมิข้อเสีย ในกรณีที่มีการบุกรุกแล้วทำการเปลี่ยนแปลงข้อมูลดังกล่าว ทำให้การตรวจจับทำไม่ได้ ดังนั้นจึงควรเก็บข้อมูลดังกล่าวไว้ในที่ๆ ปลอดภัยด้วย

2.3.1.2 การเก็บข้อมูลของการทำงานของเครื่อง

สำหรับการเก็บข้อมูลของการทำงานของเครื่อง จะเน้นไปในการเก็บข้อมูลของระบบปฏิบัติการเป็นหลัก ข้อมูลที่เก็บได้จะอยู่ในรูปของการแจ้งเตือนในระบบเช่นการตั้งค่าบางอย่างไม่สมบูรณ์ การทำงานของโปรแกรมบางโปรแกรมมีปัญหาหรือปัญหาของฮาร์ดแวร์เป็นต้น หรืออาจอยู่ในรูปข้อมูลของการทำงานโดยปกติของระบบปฏิบัติการนั้นๆ เช่น ข้อมูลการใช้งานของยูสเซอร์แต่ละคน ใคร ทำอะไร เมื่อเวลาเท่าไร ซึ่งเมื่อนำข้อมูลเหล่านี้ไปวิเคราะห์แล้ว จะได้ผลการวิเคราะห์ในลักษณะมีการใช้งานอย่างไม่ถูกต้องหรือไม่ ถ้ามี ใครเป็นผู้ใช้งานนั้นๆ เมื่อเวลาเท่าไรจากที่ไหน ข้อดีอีกข้อหนึ่งก็คือการเก็บข้อมูลในลักษณะนี้สามารถเก็บข้อมูลที่ถูกต้องเข้ารหัสได้ด้วย ส่วนข้อเสียของการเก็บข้อมูลการทำงานของเครื่องก็คือ ข้อมูลที่ได้มักจะมีขนาดใหญ่ ระบบที่ทำการเก็บข้อมูลลักษณะนี้จะมี Overhead สูงขึ้น นอกจากนี้โปรแกรมที่ทำการเก็บข้อมูลและวิเคราะห์ข้อมูลยังขึ้นอยู่กับ platform และมีราคาสูงมากด้วย

2.3.1.3 การเก็บข้อมูลการเปลี่ยนแปลงข้อมูลในระบบ

การเก็บข้อมูลการเปลี่ยนแปลงข้อมูลในระบบจะใช้หลักการของ integrity analysis ในการตรวจสอบการเปลี่ยนแปลงข้อมูลต่างๆ ในระบบ ในระบบตรวจจับการบุกรุกบางระบบจะใช้ checksum เป็นตัวบ่งบอกการเปลี่ยนแปลงในระบบ การวิเคราะห์ลักษณะนี้จะเริ่มจากการสร้างฐานข้อมูล signature ของไฟล์ต่างๆ ในระบบปกติไว้ ในระบบปกติไว้เมื่อระบบมีการทำงานก็จะทำการตรวจสอบค่า signature นี้ไปเรื่อยๆ อาจเป็นวันละครั้ง สองครั้ง หรือบ่อยกว่านั้นแล้วแต่ความสำคัญของระบบ เมื่อมีข้อมูลไฟล์ไหนมีการเปลี่ยนแปลงก็จะทราบได้ ข้อดีของการทำ

integrity analysis ลักษณะนี้ช่วยให้การตรวจจับการบุกรุกที่มีการเปลี่ยนแปลงระบบ เช่น ทำการเจาะระบบแล้วทำการวางโทรจัน หรือ Back Door ไว้ได้ สำหรับการแก้ไขเมื่อทราบว่ามีการบุกรุกมาเปลี่ยนแปลงไฟล์ข้อมูลในระบบก็ทำการแก้ไขเฉพาะไฟล์ที่ถูกแก้ไขเท่านั้น ไม่จำเป็นต้องทำการติดตั้งระบบใหม่ แต่ระบบนี้ก็ยังมีข้อเสียในกรณีที่เมื่อมีไฟล์ในระบบเยอะมาก การเก็บข้อมูล signature ของไฟล์ต่างๆ และการวิเคราะห์ข้อมูลก็จะใช้เวลานาน ระบบนี้จึงไม่เหมาะในการทำงาน real time เพราะทำให้เกิด overhead ในระบบสูงมาก

2.3.1.4 การเก็บข้อมูลเครือข่าย

การเก็บข้อมูลเครือข่ายนั้นนับวันจะมีความสำคัญขึ้นเรื่อยๆ เพราะการบุกรุกทางเครือข่ายมีมากขึ้นเรื่อยๆ การเก็บข้อมูลแบบนี้จะใช้การดักจับข้อมูลที่ผ่านไปมาในเครือข่าย โดยการทำให้เน็ตเวิร์คการ์ดอยู่ใน promiscuous mode เมื่อเน็ตเวิร์คการ์ดอยู่ในโหมดดังกล่าว จะสามารถรับข้อมูลทุกอย่างที่อยู่ในเครือข่ายได้ การเก็บข้อมูลเครือข่ายในลักษณะนี้สามารถตรวจจับการโจมตีทางเครือข่ายได้ เช่น การทำ SYN flood การทำ port scan หรือการส่งแพ็กเก็ตปริมาณมากมารบกวนในระบบ แต่เนื่องจากการเก็บข้อมูลเครือข่ายนี้ใช้ลักษณะการนำสนิฟเฟอร์เป็นหลัก จึงไม่สามารถทำงานในเครือข่ายที่เป็นเครือข่ายสวิตซ์ซึ่งไม่สามารถทำงานในระบบเครือข่ายที่เข้ารหัสข้อมูล หรือไม่สามารถเก็บข้อมูลในเครือข่ายที่มีข้อมูลหนาแน่น ได้เพราะการทำงานในการเก็บข้อมูลอาจไม่เร็วพอที่จะเก็บข้อมูลทั้งหมดที่ผ่านไปมาในระบบได้ ข้อเสียอีกข้อหนึ่งของการเก็บข้อมูลนี้ก็คือข้อมูลที่เก็บมีขนาดใหญ่มาก โดยเฉพาะอย่างยิ่งในระบบเครือข่ายที่มีการรับส่งแพ็กเก็ตปริมาณมากอยู่ตลอดเวลา

นอกจากนี้ยังมีการเก็บข้อมูลโดยทำการเก็บข้อมูลทั้ง Application-based, Host-based และ Network-based ร่วมกันด้วย เพื่อให้ได้ข้อมูลระบบอย่างครบถ้วน และใช้ข้อมูลจากทั้งสามแหล่งมาประกอบกันในการวิเคราะห์ความผิดปกติที่เกิดขึ้นในระบบด้วย การเก็บข้อมูลในลักษณะนี้เรารวมเรียกว่า “Integrated-based”

2.3.2 การวิเคราะห์ข้อมูลระบบ

เมื่อได้ข้อมูลของระบบที่จำเป็นแล้ว ในขั้นตอนต่อมาเราก็จะนำเอาข้อมูลที่ได้อามาวิเคราะห์ว่าระบบของเรามีความผิดปกติเกิดขึ้นหรือไม่ การวิเคราะห์ข้อมูลเราสามารถแบ่งการทำงานตามรูปแบบการวิเคราะห์ข้อมูลได้ 2 รูปแบบ คือ ทำการวิเคราะห์ในขณะที่เก็บข้อมูล (Real Time) หรือ จะเก็บข้อมูลทั้งหมดไว้ก่อน แล้วจึงวิเคราะห์ข้อมูลนั้นๆ ภายหลัง (Batch) ในการวิเคราะห์ข้อมูลทั้งสองรูปแบบก็มีข้อเสียแตกต่างกันไป

2.3.2.1 การวิเคราะห์ในขณะที่เก็บข้อมูล (Real Time)

ในการวิเคราะห์ข้อมูลที่ได้ในขณะที่เก็บข้อมูล หรือ แบบ Real Time นั้น ระบบจะจัดเก็บข้อมูล วิเคราะห์ข้อมูล และรายงานผลการวิเคราะห์ ในช่วงเวลาเดียวกัน เมื่อเกิดข้อผิดพลาดขึ้น

สามารถตอบสนองได้ทันทีที่ระบบที่ทำงานแบบ Real Time มีการแจ้งเตือนหลายๆ แบบ เช่น E-mail หรือ Instant Messaging ให้กับผู้ดูแลระบบได้ในช่วงเวลาที่มีการบุกรุกได้ ในการตรวจสอบระบบแบบ Real Time ทำให้ระบบสามารถตรวจสอบข้อผิดพลาดได้อย่างรวดเร็ว แต่ก็ขึ้นอยู่กับความเร็วในการวิเคราะห์ข้อมูลด้วย ถ้าข้อมูลมีความซับซ้อนมากๆ ก็จะใช้เวลามากตาม เช่นเดียวกันเมื่อระบบทำการตรวจสอบการบุกรุกได้ในขณะที่เพิ่งเกิดการบุกรุกขึ้น ผู้ดูแลระบบ หรือ ระบบตรวจจับการบุกรุกเองสามารถแก้ไขปัญหาที่เกิดขึ้นได้ทันที แต่ทั้งนี้ก็ขึ้นอยู่กับความเร็วในการวิเคราะห์ข้อมูล และชนิดของปัญหาว่ามีความยุ่งยากในการแก้ปัญหาเล็กน้อยเพียงไรด้วย

ระบบ Real Time ทำงานได้อย่างรวดเร็ว แต่การทำงานที่รวดเร็วดังกล่าวก็ต้องแลกกับการใช้หน่วยความจำปริมาณมาก และการประมวลผลที่รวดเร็วมากด้วย อีกทั้งการตอบสนองต่อการบุกรุกโดยอัตโนมัติ อาจทำให้เกิดความเสียหายกับระบบมากกว่าเดิม เพราะว่าในบางครั้ง การทำงานที่เร็วเกินไปของระบบนี้ ทำให้เกิดความผิดพลาดในการวิเคราะห์ข้อมูลจนประมวลผลการทำงานปกติ กลายเป็นการทำงานที่ผิดปกติ แล้วทำการแก้ไขตามข้อมูลที่มีอยู่ ก็ยิ่งทำให้ระบบมีความเสียหายมากกว่าเดิม ระบบตรวจจับการบุกรุกที่ทำงานแบบ Real Time จึงเหมาะกับระบบที่มีข้อมูลที่ต้องพิจารณาน้อย ต้องการการรายงานอย่างรวดเร็วเมื่อผิดปกติ และข้อมูลที่ต้องนำมาวิเคราะห์ไม่ซับซ้อนมากนัก

2.3.2.2 การวิเคราะห์ข้อมูลภายหลังจากที่เก็บข้อมูล (Batch)

อีกรูปแบบหนึ่งในการวิเคราะห์ระบบที่ใช้กัน คือ การวิเคราะห์ข้อมูลภายหลังจากที่เก็บข้อมูลไว้แล้วหรือการทำงานแบบ Batch การทำงานในแบบนี้เหมาะกับงานที่ไม่จำเป็นต้องตอบสนองทันทีเมื่อเกิดความผิดปกติขึ้น แต่ให้มีการบันทึก และรายงานว่าเกิดความผิดปกติขึ้น การทำงานจะใช้หน่วยความจำ และการประมวลผลน้อยกว่าแบบแรก แต่ก็ใช้เนื้อที่ในการเก็บข้อมูลมากกว่าแบบแรกแน่นอน เหมาะกับองค์กรที่มีบุคลากรจำกัด ข้อเสียของการทำงานแบบ Batch คือ มักแก้ปัญหาที่เกิดขึ้นไม่ทัน เพราะกว่าจะทราบว่าเกิดปัญหาขึ้น ปัญหานั้นก็เกิดขึ้นนานมาก ความเสียหายที่เกิดขึ้นก็แก้ไขได้ยาก

ไม่ว่าจะเป็นการวิเคราะห์ระบบแบบ Real Time หรือ Batch ก็จะมีวิธีการวิเคราะห์ระบบที่เหมือนกัน คือ ทำการหารูปแบบของการโจมตีในข้อมูลที่ได้รับมา (Signature Analysis) วิธีการวิเคราะห์ทางสถิติ (Statistical Analysis) และวิธีการตรวจสอบการเปลี่ยนแปลงของระบบ (Integrity Analysis)

1) การหารูปแบบของการโจมตี (Signature Analysis)

ในวิธีการวิเคราะห์ระบบแบบ Signature Analysis เป็นการวิเคราะห์ข้อมูลโดยการหาสัญญาณของการโจมตี (Attack Signature) การทำงานจะทำโดยการเปรียบเทียบรูปแบบของข้อมูลกับรูปแบบของการโจมตีในฐานข้อมูล ว่ามีความคล้ายกันหรือไม่ ถ้ามีความคล้ายคลึงกันก็แสดง

ว่ามีการโจมตีเกิดขึ้นแล้ว ในการเปรียบเทียบอาจเป็นแบบอย่างง่าย คือ การเปรียบเทียบข้อมูลว่ามีความเข้ากันได้กับข้อมูลของการโจมตีเพียงใด หรือเป็นแบบที่มีความสลับซับซ้อนขึ้นอีก เช่น การทำ state transition เป็นต้น

สำหรับโปรแกรมตรวจจับการบุกรุกที่มีจำหน่ายในท้องตลาด ส่วนใหญ่จะทำงานในลักษณะของการเปรียบเทียบรูปแบบกับการโจมตีในฐานข้อมูล ซึ่งบริษัทผู้ขายจะให้ฐานข้อมูลของการโจมตีไว้ด้วย ผู้ใช้งานจะมีการอัปเดตข้อมูลในฐานข้อมูลบ่อยๆ เพื่อเพิ่มความสามารถในการวิเคราะห์ข้อมูลในระบบ การวิเคราะห์ระบบด้วยวิธี Signature analysis นี้จะมี overhead ไม่มากนักเพราะเป็นเพียงการเปรียบเทียบข้อมูลกับข้อมูลในฐานข้อมูลเท่านั้น และยังเพิ่มความเร็วในการทำงานโดยรวมมากขึ้น เพราะสามารถนำข้อมูลในฐานข้อมูลเป็นกฎในการกรองข้อมูลที่จะเก็บให้น้อยลงไปด้วย แต่วิธีการนี้ก็มีข้อเสียเพราะฐานข้อมูลจะมีขนาดใหญ่ขึ้นเรื่อยๆ ต้องมีการอัปเดตฐานข้อมูลบ่อยๆ

2) วิธีการวิเคราะห์ทางสถิติ (Statistical Analysis)

วิธีการวิเคราะห์ทางสถิติ (Statistical Analysis) เป็นวิธีการวิเคราะห์ข้อมูลอีกแบบหนึ่งที่มีแนวคิดตรงข้ามกับวิธีการแรก คือ จะหารูปแบบของการทำงานที่เป็นปกติ แล้วสร้างเป็นโพรไฟล์ (Profile) เก็บไว้ก่อน ในการวิเคราะห์ข้อมูล จะเปรียบเทียบข้อมูลกับโพรไฟล์ที่สร้างไว้ ถ้าไม่เข้ากันก็แสดงว่ามีความผิดปกติเกิดขึ้นแล้ว สำหรับโพรไฟล์นั้นอาจแยกเป็นเป็นโพรไฟล์สำหรับแอปเจ็คต์ต่างๆ ในระบบ เช่น ยูสเซอร์ ไฟล์ ไคลเอนต์ และอุปกรณ์ต่างๆ โดยรวมละเอียดที่เก็บอยู่ในโพรไฟล์จะมีข้อมูลของจำนวนครั้งที่เข้าสู่ระบบ จำนวนครั้งที่เข้าสู่ระบบผิดพลาด เวลา และข้อมูลอื่นๆ ที่จำเป็น ค่าแต่ละค่าที่เก็บจะเป็นค่าของการใช้งานที่เป็นปกติ การตรวจจับว่าเกิดความผิดปกติขึ้นแล้ว จะดูจากค่าที่ไม่เข้ากับโพรไฟล์ เป็นต้น ยกตัวอย่างเช่น ในการทำงานปกติ ผู้ใช้งานฐานข้อมูลจะมีการเข้าใช้ข้อมูลในฐานข้อมูลตั้งแต่เวลา 8 โมงเช้า ถึง 6 โมงเย็นเท่านั้น แต่ข้อมูลที่ตรวจจับได้มีการเข้าใช้ฐานข้อมูลตอนดึกสอง ซึ่งก็สามารถบอกได้ว่าการบุกรุกเกิดขึ้นแล้ว

เนื่องจากวิธีการวิเคราะห์ข้อมูลแบบ Statistical Analysis เป็นการตรวจจับโดยใช้หลักการของการอนุญาตให้ใช้งานในการทำงานทุกๆ ไปและไม่อนุญาตให้ใช้งานนอกเหนือจากที่เคยใช้ โดยทั่วไปเท่านั้น การตรวจจับในลักษณะนี้จึงสามารถตรวจจับการบุกรุกที่ไม่เคยเจอมาก่อนได้ และสามารถตรวจจับการบุกรุกในรูปแบบที่ซับซ้อนได้ด้วย เพราะเรารู้ว่าการทำงานที่สลับซับซ้อนมักจะไม่เหมือนกับการทำงานโดยปกติ แต่วิธีการวิเคราะห์ข้อมูลแบบนี้ก็มีข้อเสีย เนื่องจากการเก็บข้อมูลการทำงานที่เป็นปกติไว้เพื่อเปรียบเทียบกับการทำงานที่ผิดปกติ เมื่อทำการบุกรุกในลักษณะเดิมๆ เป็นเวลานาน ก็จะทำให้โพรไฟล์มีการเปลี่ยนแปลง ระบบตรวจจับก็จะเห็นว่าการโจมตีในลักษณะนั้นเป็นการทำงานที่เป็นปกติแทน และไม่สามารถตรวจจับการทำงานที่ผิดปกติในลักษณะนั้นได้อีกต่อไป และไม่เหมาะกับองค์กรที่มีการเปลี่ยนแปลงการ

ทำงานบ่อยๆ เพราะทำให้โพรไฟล์มีขนาดใหญ่ ทำให้ระบบตรวจจับรวน และมีความผิดพลาดในการวิเคราะห์สูง

3) วิธีการตรวจสอบการเปลี่ยนแปลงของระบบ (Integrity Analysis)

วิธีสุดท้ายที่นิยมใช้กันในการวิเคราะห์ข้อมูลคือ วิธีการตรวจสอบการเปลี่ยนแปลงของระบบ (Integrity Analysis) ลักษณะการทำงานของวิธีการนี้ คือ การหาว่ามีการเปลี่ยนแปลงเกิดขึ้นในระบบหรือไม่ เช่น มีไฟล์ไหนมีการเปลี่ยนแปลง หรือมีออปเจกต์อะไรที่มีการเปลี่ยนแปลงคุณสมบัติบ้าง แล้วทำการแจ้งเตือนกับผู้ดูแลระบบ ในการวิเคราะห์ลักษณะนี้จะใช้แฮชอัลกอริทึม (hash algorithm) เพื่อสร้างเมสเสจไดเจส (message digest) ของข้อมูล แล้วทำการเปรียบเทียบเมสเสจไดเจส ของข้อมูลในช่วงเวลาต่างๆ ว่าเหมือน หรือต่างกันหรือไม่ ถ้าเมสเสจไดเจสต่างก็แสดงว่าข้อมูลมีการเปลี่ยนแปลง Integrity Analysis สามารถตรวจจับการบุกรุกที่เข้ามาเปลี่ยนแปลงข้อมูลในระบบ หรือมีการติดตั้ง โปรแกรม เช่น sniffer rootkit ต่างๆ ในระบบได้ แต่ก็มีข้อเสีย คือ การวิเคราะห์ระบบในลักษณะนี้จะทำงานเป็นแบบ batch เท่านั้น ไม่เหมาะกับการทำ real time อย่างยิ่งเพราะจะทำให้เปลืองทรัพยากรมาก

2.3.3 การตอบสนอง

เมื่อมีการตรวจพบว่าการบุกรุกเกิดขึ้นในระบบ สำหรับระบบตรวจจับที่ทำงานแบบ real time จะมีการตอบสนองต่อการบุกรุกเพื่อ ไม่ให้เกิดความเสียหาย หรือบรรเทาความเสียหายที่เกิดขึ้นสำหรับระบบที่ทำงานเป็นแบบ batch การตอบสนองอาจทำได้ไม่มากนัก เพราะการบุกรุกนั้นเกิดไปแล้ว ความเสียหายก็เกิดขึ้นแล้ว การตอบสนองอาจอยู่ในรูปแบบการบรรเทาไม่ให้ความเสียหายมีเพิ่มมากขึ้นเท่านั้น การตอบสนองต่อการบุกรุกนั้นแบ่งออกได้เป็นสามแบบด้วยกัน คือ การเปลี่ยนแปลงสภาพของระบบ การแก้ไขความผิดพลาดให้ถูก และการแจ้งเตือนผู้ดูแลระบบเมื่อถูกบุกรุก

2.3.3.1 การเปลี่ยนแปลงสภาพของระบบ

สำหรับการตอบสนองต่อการบุกรุกโดยการเปลี่ยนแปลงสภาพของระบบที่ถูกโจมตีก็เพื่อแก้ปัญหาหรือลดความเสียหายที่จะเกิดขึ้น เช่น ตัดการเชื่อมต่อระหว่างระบบกับการบุกรุกออกจากกัน การตั้งค่าอุปกรณ์เครือข่าย หรือไฟร์วอลล์ไม่ให้มีการติดต่อกับระบบของการบุกรุกอีกต่อไป และการหาข้อมูลเกี่ยวกับการโจมตีโดยอัตโนมัติเพื่อตรวจหาการบุกรุกต่อไป

2.3.3.2 การแก้ไขความผิดพลาดให้ถูก

การแก้ไขระบบ เป็นการตอบสนองต่อปัญหาที่เกิดขึ้นแล้วในระบบ โดยปกติแล้วการบุกรุกมักเปลี่ยนแปลงค่าต่างๆ ในระบบ โดยเฉพาะเข้ามาทำการเปลี่ยนแปลงข้อมูลในระบบตรวจจับการบุกรุกเพื่อไม่ให้สามารถตรวจจับการบุกรุกได้ การแก้ไขระบบก็เพื่อให้ระบบดังกล่าวสามารถทำงานได้อย่างเป็นปกติ

2.3.3.3 การแจ้งเตือนผู้ดูแลระบบ

สุดท้ายเป็นการแจ้งเตือนผู้ดูแลระบบ โดยปกติมักแจ้งเตือนผู้ดูแลทันทีเมื่อทำการวิเคราะห์ได้มีความผิดปกติเกิดขึ้น เพื่อให้ผู้ดูแลระบบรับรู้และสามารถแก้ไขระบบได้ทันทีที่สำหรับการแจ้งเตือนนี้ ผู้ดูแลระบบสามารถเลือกได้ว่าจะแจ้งเตือนใครบ้าง และทำการแจ้งเตือนในรูปแบบไหนอาจเป็น E-mail, Pager หรือ Instant Messaging ต่างๆ

2.3.4 การรายงานผลการทำงาน

เมื่อระบบตรวจจับการบุกรุกทำการวิเคราะห์ระบบ และตรวจพบความผิดปกติในระบบ อาจมีการตอบสนองต่อความผิดปกตินั้นถ้าทำได้ จากนั้นระบบตรวจจับการบุกรุกต้องมีการรายงานผลให้กับผู้ดูแลระบบทราบในรูปแบบต่างๆ โดยรายละเอียดของการรายงานผลนั้น จะบอกถึงช่องโหว่ในระบบ การแก้ไขปัญหาคว่าๆ บางครั้งอาจมีรายละเอียดของความรู้พื้นฐานบางอย่างของระบบ ที่ทำให้เกิดการบุกรุกลักษณะนั้นๆ ได้ การรายงานผลการทำงาน นอกจากเป็นการรายงานต่อผู้ดูแลระบบเพื่อให้ทราบการทำงาน หรือจุดอ่อนในระบบแล้ว ยังเป็นประโยชน์ต่อการวิเคราะห์สถานะของระบบ และการวิเคราะห์ความปลอดภัยในระบบอีกด้วย

2.4 ความสำคัญของระบบตรวจจับการบุกรุก

เมื่อได้ทราบถึงการทำงานคร่าวๆ ของระบบตรวจจับการบุกรุกแล้ว อาจคิดว่าระบบตรวจจับการบุกรุกไม่มีความสำคัญเพราะในเมื่อมีการใช้งานไฟร์วอลล์อยู่แล้ว แต่ความเป็นจริงแล้วถึงแม้ว่าระบบจะมีไฟร์วอลล์อยู่แล้วก็ยังจำเป็นต้องใช้ระบบตรวจจับการบุกรุกด้วยเพราะในงานบางอย่างไฟร์วอลล์ก็ไม่สามารถช่วยได้

จุดประสงค์ของการใช้งานไฟร์วอลล์นั้น สร้างขึ้นเพื่อเป็นเสมือนตัวป้องกันระบบให้แยกตัวออกมาจากเครือข่ายที่ไม่ปลอดภัย เป็นเหมือนเมืองหน้าด่านของระบบ เป็นผู้ป้องกันการบุกรุกจากภายนอก แต่ระบบตรวจจับการบุกรุกนั้นมีจุดประสงค์ที่แตกต่างไป โดยเป็นผู้เฝ้าดูระบบ และเป็นผู้เตือนเมื่อเกิดความผิดปกติเกิดขึ้น ยกตัวอย่างในอาคารใหญ่ๆ จะมี คนคอยดูแลอยู่ภายนอกกับคนที่ไม่ควรเข้ามาในอาคารให้อยู่ภายนอก แต่ภายในตัวอาคารก็จะมีกล้องวิดีโอคอยตรวจตราอยู่ภายในมีกริ่งสัญญาณเตือนเมื่อเกิดความผิดปกติเกิดขึ้น ซึ่งก็ช่วยให้แก้ปัญหาได้ทันที และในกรณีที่มีความผิดปกติเกิดขึ้น แต่ไม่สามารถตรวจจับได้ในขณะนั้น ระบบตรวจจับการบุกรุกก็มีการจัดเก็บข้อมูลการใช้ระบบไว้ จึงสามารถนำข้อมูลดังกล่าวมาวิเคราะห์หาความผิดปกติได้ภายหลัง โดยจุดประสงค์ในการสร้างระบบตรวจจับการบุกรุกและไฟร์วอลล์ จึงต่างกันโดยสิ้นเชิง แต่ถึงแม้ว่าจุดประสงค์การทำงานของระบบตรวจจับการบุกรุกและไฟร์วอลล์ จะแตกต่างกัน แต่ทั้งสองก็สามารถทำงานร่วมกันและทำให้ประสิทธิภาพการรักษาความปลอดภัยในระบบดีขึ้นด้วย

2.5 สรุป

ระบบตรวจจัดการบุงรุกเป็นผู้ช่วยที่ดีสำหรับผู้ดูแลระบบ หน้าที่ของการตรวจจัดการบุงรุกนั้นจะรวบรวมข้อมูล วิเคราะห์ข้อมูลและทำการแจ้งผลการวิเคราะห์ข้อมูลต่างๆ ในระบบให้กับผู้ดูแลระบบ ซึ่งช่วยให้การดูแลระบบทำได้อย่างมีประสิทธิภาพ แต่ระบบตรวจจัดการบุงรุกก็ยังไม่ใช้สิ่งที่จะมาแก้ปัญหาคือความปลอดภัยในระบบ โดยสิ้นเชิงได้ เนื่องจากการที่จะทำให้ระบบปลอดภัยต้องอาศัยความร่วมมือจากหลายๆ ฝ่าย ไม่เพียงแต่เฉพาะการดูแลของผู้ดูแลระบบคนเดียว สิ่งที่สำคัญที่สุด ที่ผู้เขียนเห็นว่าจะสร้างความปลอดภัยให้กับระบบในระยะยาวก็คือ การปลูกจิตสำนึกด้านความปลอดภัยในการใช้งานให้กับผู้ใช้ระบบแต่ละคนด้วย

ส่วนในงานวิจัยนี้เป็นการประยุกต์ระบบตรวจจัดการบุงรุกแบบ Anomaly โดยใช้ SOM ซึ่งกล่าวรายละเอียดในบทที่ 4 ต่อไป

บทที่ 3

นิเวศเน็ตเวิร์กแบบไม่มีผู้สอน

ในบทนี้จะกล่าวถึงโมเดลนิเวศเน็ตเวิร์กของโคโฮเนน (Kohonen) ซึ่งเป็นนิเวศเน็ตเวิร์กแบบไม่มีผู้สอนที่ได้รับความนิยมเป็นอย่างมาก โดยจะอธิบายถึงการทำงานของเซลล์ในเน็ตเวิร์กเมื่อได้รับอินพุตเข้ามา ชั้นตอนต่าง ๆ ในการเรียนรู้รวมทั้งแสดงตัวอย่างงานวิจัยที่นำเอาโมเดลนี้ไปประยุกต์ใช้ในด้านต่าง ๆ

3.1 บทนำ

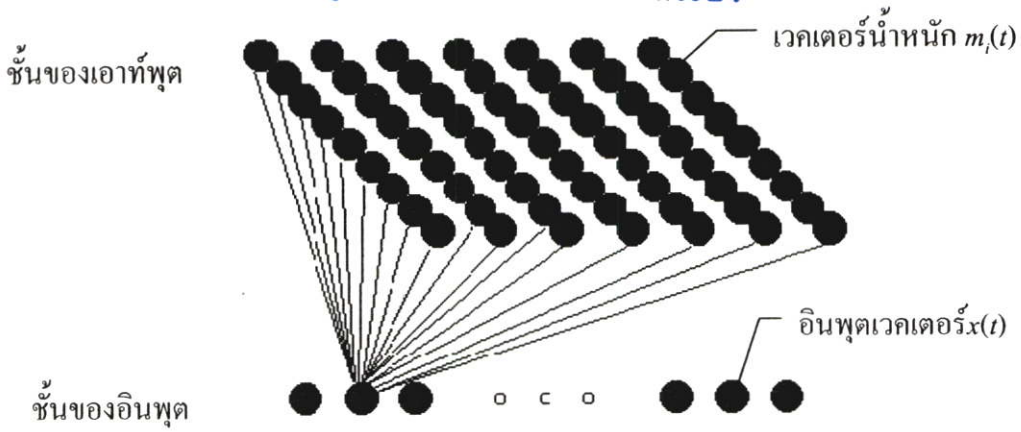
เซลฟออร์แกนไนซิงแมป (Self-Organizing Map) หรือ SOM เป็นนิเวศเน็ตเวิร์กแบบไม่มีผู้สอนประเภทหนึ่ง [8] ซึ่งแตกต่างจากนิเวศเน็ตเวิร์กแบบมัลติเลเยอร์เพอเซพตอล (Multi Layer Perceptron MLP) หรือแบบแบ็คพรอพาเกชัน (Backpropagation) ซึ่งเป็นนิเวศเน็ตเวิร์กแบบมีผู้สอน อัลกอริทึมของ SOM ถูกนำเสนอโดยศาสตราจารย์โคโฮเนน (Kohonen) ในปี ค.ศ. 1982 ซึ่งรู้จักกันในชื่อ แผนภาพคุณลักษณะของโคโฮเนน (Kohonen feature map)

SOM เป็นอัลกอริทึมที่ใช้ในการจัดกลุ่มข้อมูล โดยสามารถจัดกลุ่มข้อมูลที่มีมิติสูงให้อยู่ในรูปของแผนภาพ 2 มิติซึ่งประกอบไปด้วยโหนดของนิเวศเน็ตเวิร์ก ข้อมูลจะถูกจัดลงในโหนดต่าง ๆ ของแผนภาพ หลังจากเสร็จสิ้นกระบวนการเรียนรู้แผนภาพจะถูกจัดเรียงตัว โดยข้อมูลที่มีความคล้ายคลึงกันจะอยู่ในกลุ่มโหนดใกล้เคียงกัน ดังนั้นแผนภาพที่ได้จะแสดงคุณลักษณะของข้อมูลได้เป็นอย่างดี เช่น การกระจายของข้อมูล ความสัมพันธ์ระหว่างข้อมูลในกลุ่ม เป็นต้น

3.2 โมเดลและอัลกอริทึมของ SOM

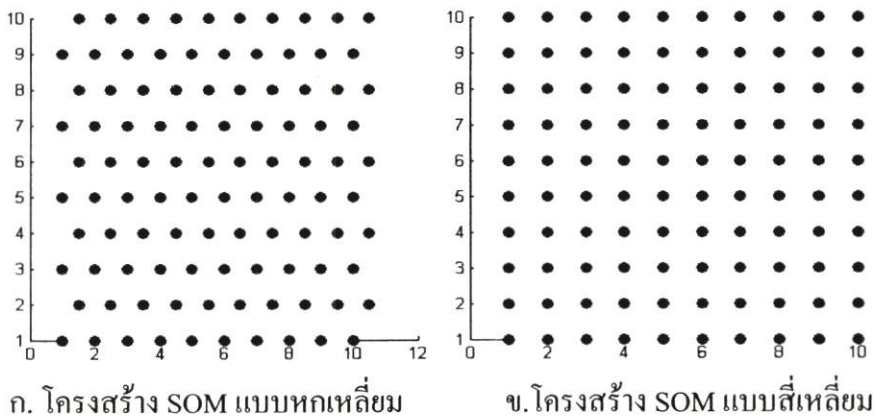
โมเดลของ SOM ประกอบด้วยเซลล์ 2 ชั้น [9] ดังรูปที่ 3.1 ชั้นแรกคือชั้นของอินพุต (Input layer) ประกอบด้วยเซตของอินพุตเวกเตอร์ $x(t)$ ที่มีขนาด n มิติ ($1 \times n$ มิติ) ซึ่งเป็นอินพุตที่ใช้ในการเรียนรู้ของแผนภาพ โดยที่ t คืออินเด็กซ์ของอินพุตหรือแทน เวลาใด ๆ ก็ได้ ชั้นที่สองคือชั้นของแผนภาพโคโฮเนน (Kohonen layer) หรือชั้นของเอาต์พุตประกอบด้วยโหนดของนิเวศเน็ตเวิร์กที่เรียงตัวอยู่ในรูปแบบของแผนภาพ 2 มิติ ในแต่ละโหนด i จะเป็นค่าเวกเตอร์นำหน้าแทนด้วย $m_i(t)$ นั่นคือ $m_i(t) \in \mathcal{R}^n$ โดยที่ \mathcal{R}^n คือ โดเมนของขนาดของ n และขนาดของเวกเตอร์นำหน้าจะต้องมีขนาดเท่ากับอินพุตเวกเตอร์ $x(t)$

สำนักหอสมุดกลาง พระจอมเกล้าลาดกระบัง



รูปที่ 3.1 แสดงโมเดลพื้นฐานของ SOM แบบสี่เหลี่ยม

ในการออกแบบโครงสร้างโมเดล SOM เราสามารถกำหนดโครงสร้างของโหนดได้ดังรูปที่ 3.2 ก. เป็นการกำหนดโครงสร้างของ SOM แบบหกเหลี่ยม ในรูปที่ 3.2 ข. เป็นการกำหนดโครงสร้างของ SOM แบบสี่เหลี่ยม ซึ่งทั้งสองจะมีการกำหนดโหนดใกล้เคียงที่ต่างกัน โดยที่โครงสร้างแบบหกเหลี่ยมจะมีโหนดใกล้เคียงเป็นรูปหกเหลี่ยมมี แต่โครงสร้างแบบสี่เหลี่ยมจะมีโหนดใกล้เคียงเป็นรูปสี่เหลี่ยม



ก. โครงสร้าง SOM แบบหกเหลี่ยม

ข. โครงสร้าง SOM แบบสี่เหลี่ยม

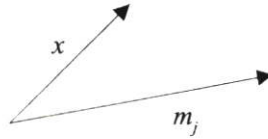
รูปที่ 3.2 แสดงโครงสร้างของ SOM

กระบวนการเรียนรู้ของ SOM เกิดขึ้นจากการปรับตัวของเวกเตอร์น้ำหนักที่มีต่ออินพุตเวกเตอร์ โดยเริ่มแรกจะทำกำหนดน้ำหนักเริ่มต้นขนาดเล็กให้กับโหนดทุกโหนด จากนั้นจะเริ่มค้นกระบวนการเรียนรู้ดังนี้

1. เลือกอินพุตเวกเตอร์แบบสุ่มเลือกจากอินพุตโดเมน
2. เปรียบเทียบอินพุตเวกเตอร์ $x(t)$ กับ โหนด $m_i(t)$ ทุกโหนดเพื่อหาโหนดชนะจากโหนดทั้งหมด
3. ปรับเวกเตอร์น้ำหนักของโหนดชนะ เพื่อให้โหนดชนะเข้าใกล้อินพุตมากขึ้น

4. ปรับเวกเตอร์น้ำหนักของโหนดใกล้เคียง เพื่อให้อินพุตเวกเตอร์ถัดไปที่มีค่าใกล้เคียง มีโหนดชนะใหม่อยู่ใกล้กัน

กระบวนการเหล่านี้จะถูกทำซ้ำไปเรื่อย ๆ จนกว่าจะสอดคล้องตามเงื่อนไขหรือจนกว่าจะครบจำนวนรอบของการเรียนรู้ จากกระบวนการเรียนรู้ข้างต้นมีการคำนวณที่สำคัญอยู่ 2 ส่วนคือ ส่วนแรกคือการคำนวณเพื่อหาโหนดชนะ(ขั้นตอนที่ 2) ในการคำนวณหาโหนดชนะอินพุตเวกเตอร์ $x(t)$ ถูกนำไปเปรียบเทียบกับโหนด $m_j(t)$ ทุกโหนดเพื่อหาโหนดชนะจากโหนดทั้งหมด ฟังก์ชันที่ใช้ในการเปรียบเทียบโดยทั่วไปแล้วจะใช้ฟังก์ชันวัดระยะทางแบบยูคลิด (Euclidean distance) ดังรูปที่ 3.3



รูปที่ 3.3 แสดงระยะทางแบบยูคลิดระหว่างเวกเตอร์ x และ m_j

การหาโหนดที่ชนะ c สามารถหาได้จากโหนดที่มีระยะห่างระหว่างอินพุตเวกเตอร์กับเวกเตอร์น้ำหนักของโหนดนั้นน้อยที่สุดดังสมการที่ 1.1

$$c : m_c(t) = \min_i \| x(t) - m_i(t) \| \quad (1.1)$$

ส่วนที่สองคือการปรับเวกเตอร์น้ำหนัก หลังจากที่ได้โหนดชนะแล้วจะต้องทำการปรับน้ำหนักเพื่อให้เข้าใกล้อินพุตมากขึ้น นอกจากการเรียนรู้ที่เกิดขึ้นที่โหนดชนะแล้ว โหนดใกล้เคียง จะเกิดการเรียนรู้ด้วย ค่าเวกเตอร์น้ำหนักของโหนดใกล้เคียงจะปรับค่าให้เข้าใกล้กับอินพุตเวกเตอร์เดียวกัน เพื่อเพิ่มโอกาสให้อินพุตใหม่ที่ใกล้เคียงกับอินพุตเดิมสามารถที่จะมีโหนดชนะใหม่ใกล้กับโหนดชนะเดิมได้ สมการในการปรับค่าน้ำหนักสามารถแสดงได้ดังสมการที่ 1.2

$$m_j(t+1) = m_j(t) + \alpha(t) \times h_{c,j}(t) \times [x(t) - m_j(t)] \quad (1.2)$$

เมื่อ

- t คือรอบปัจจุบันของการเรียนรู้
- $x(t)$ คืออินพุตเวกเตอร์ปัจจุบัน
- $m_j(t)$ คือเวกเตอร์น้ำหนัก
- $\alpha(t)$ คืออัตราการเรียนรู้

โดยที่อัตราการเรียนรู้ $\alpha(t)$ จะขึ้นอยู่กับจำนวนรอบซึ่งแสดงเป็นสมการเชิงเส้น ได้ดังสมการที่ 1.3

$$\alpha(t) = \alpha(0) \times \frac{T-t}{T} \quad (1.3)$$

เมื่อ

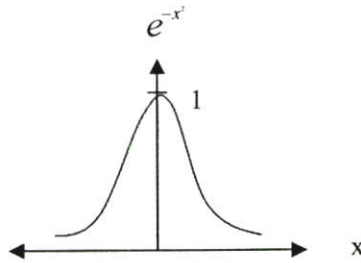
- T คือจำนวนรอบทั้งหมด
- t คือจำนวนรอบปัจจุบัน

$h_{c_i}(t)$ คือฟังก์ชันที่ใช้ในการกำหนดน้ำหนักในการปรับค่าโหนดใกล้เคียงโดยทั่วไปแล้ว $h_{c_i}(t)$ จะใช้ฟังก์ชันเกาส์เซียน (Gaussian) ซึ่งสามารถเขียนได้ดังสมการที่ 1.4

$$h_{c_i}(t) = \exp\left(-\frac{\|r_c - r_i\|^2}{2\sigma^2(t)}\right) \quad (1.4)$$

เมื่อ

- $\|r_c - r_i\|$ คือระยะห่างของตำแหน่งของโหนด i กับโหนดชนะ c
- $\sigma(t)$ คือรัศมีของบริเวณโหนดใกล้เคียง

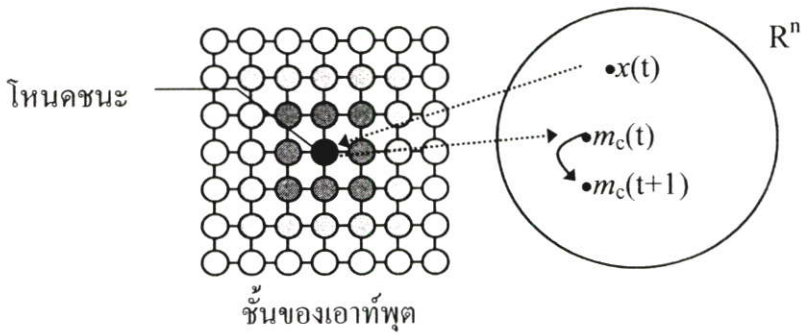


รูปที่ 3.4 แสดงกราฟของฟังก์ชัน Gaussian ($y=e^{-x^2}$)

ในรูปที่ 3.4 ลักษณะของฟังก์ชันเกาส์เซียนคือ เมื่อค่า x คือค่าระยะห่างมีค่ามาก ค่าที่ส่งกลับมาจากฟังก์ชันจะลดลงไปเรื่อย ๆ จนเข้าใกล้ศูนย์ ซึ่งสอดคล้องกับการปรับน้ำหนักของโหนดชนะและโหนดใกล้เคียง โหนดชนะจะมีค่า x เป็นศูนย์ซึ่งจะให้ค่าเกาส์เซียนฟังก์ชันออกมาเป็นหนึ่งซึ่งมากที่สุด โหนดที่ใกล้กับโหนดชนะจะมีการปรับค่าเวคเตอร์น้ำหนักมากกว่าโหนดที่อยู่ไกล โดยจะมีการกำหนดรัศมีของโหนดใกล้เคียง

โดยปกติรัศมีของโหนดใกล้เคียงจะค่อย ๆ ลดลงตามจำนวนรอบในการเรียนรู้ t ดังสมการที่ 1.5

$$\sigma(t+1) = 1 + (\sigma(t) - 1) \times \frac{T-t}{T} \quad (1.5)$$

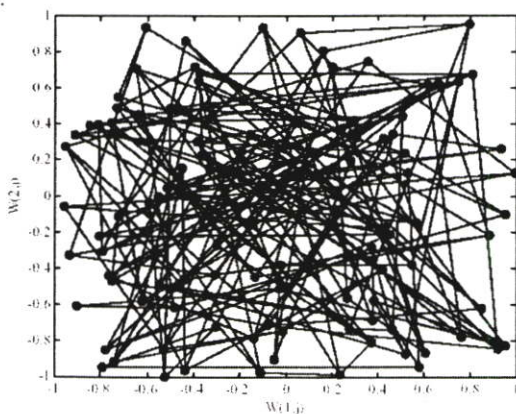


รูปที่ 3.5 แสดงโครงสร้างของ SOM ขนาด 7x7

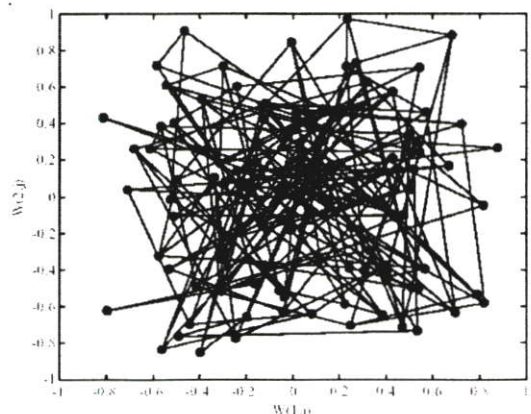
ในรูปที่ 3.5 แสดง SOM ขนาด 7x7 แบบสี่เหลี่ยมโหนดสี่เหลี่ยมที่เชื่อมที่สุดคือโหนดชนะสำหรับอินพุตเวกเตอร์ $x(t)$ จากนั้นค่าเวกเตอร์น้ำหนักของโหนด $m_c(t)$ จะถูกปรับค่าให้เข้าใกล้กับอินพุตเวกเตอร์มากขึ้น หลังจากนั้นจะทำการปรับโหนดใกล้เคียงของโหนดชนะ โดยความเข้มสีของโหนดจะแสดงถึงปริมาณการปรับค่าของเวกเตอร์น้ำหนัก โหนดที่มีสีเข้มมากจะมีการปรับค่าเวกเตอร์น้ำหนักมากกว่าสีเข้มน้อย

ตัวอย่างที่ 3.1 การใช้แผนภาพ SOM ในการเรียนรู้อินพุตเวกเตอร์ 2 มิติ

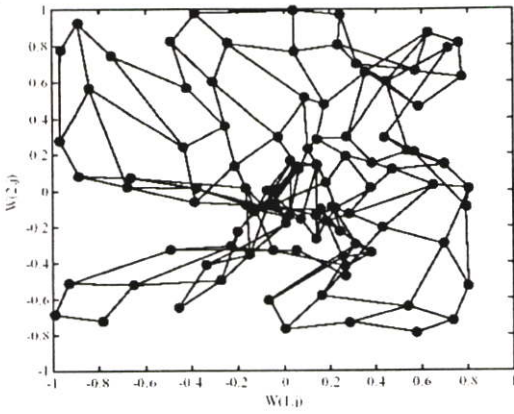
ในตัวอย่างนี้แสดงการเรียนรู้ของแผนภาพ SOM ขนาด 10x10 โดยจะทำการจัดกลุ่มอินพุตเวกเตอร์ที่มีขนาด 2 มิติ จำนวน 1000 เวกเตอร์ โดยอินพุตเวกเตอร์ได้จากการสุ่มค่าที่อยู่ในช่วงของ $[-1,1]$ และเวกเตอร์น้ำหนักของโหนดได้จากการสุ่มอยู่ในช่วงของ $[-1,1]$ ด้วย อัตราการเรียนรู้เริ่มต้น $\alpha=0.1$ ผลของการเรียนรู้แสดงในรูปที่ 3.6



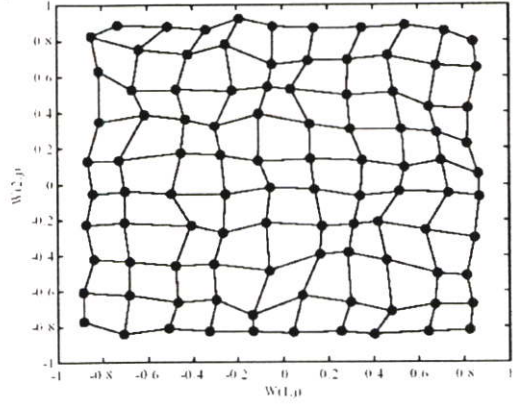
ก. แผนภาพเริ่มต้นแบบสุ่ม



ข. แผนภาพ SOM หลังจาก 100 รอบ



ค. แผนภาพ SOM หลังจาก 1000 รอบ



ง. แผนภาพ SOM หลังจาก 10000 รอบ

รูปที่ 3.6 แสดงแผนภาพ SOM ณ จำนวนรอบที่แตกต่างกัน

ในรูป 3.6 แสดงตัวอย่างแผนภาพ SOM ณ จำนวนรอบที่แตกต่างกัน จุดในรูปจะแสดงเวกเตอร์นำหนักของโหนด w_{1j} , w_{2j} ผลลัพธ์ที่ได้ในรูป 3.6 ง. เมื่อเสร็จสิ้นกระบวนการแผนภาพจะถูกจัดเรียงอย่างถูกต้อง โดยอินพุตเวกเตอร์ 1 ตัวก็จะตอบสนองกับโหนดเพียง 1 โหนด แต่โหนด 1 โหนดอาจจะตอบสนองกับอินพุตเวกเตอร์มากกว่า 1 ตัวก็ได้ หรือ บางโหนดอาจจะไม่ตอบสนองกับอินพุตเวกเตอร์ใด ๆ เลย

ตัวอย่างที่ 3.2 การใช้งานแผนภาพ SOM ในการเรียนรู้อินพุตเวกเตอร์ RGB

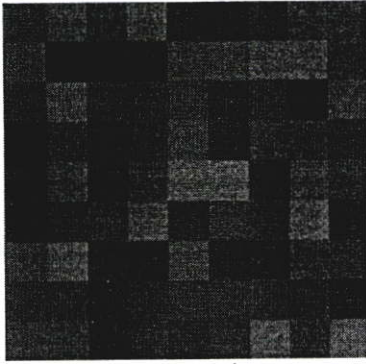
เรากำหนดแผนภาพ SOM ให้มีขนาด 9×9 เซตข้อมูลที่ใช้การเรียนรู้เป็นเวกเตอร์ RGB จำนวน 500 เวกเตอร์

ตารางที่ 3.1 แสดงอินพุตเวกเตอร์ในรูปแบบของ RGB

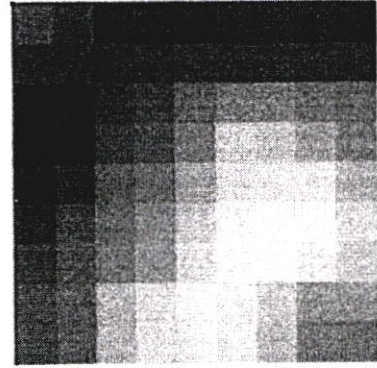
R	G	B
250	235	215
165	042	042
210	105	30
255	140	0
233	150	122
...

ตารางที่ 3.1 แสดงตัวอย่างของข้อมูลสีซึ่งสามารถเขียนเป็นอินพุตเวกเตอร์ที่มีขนาด 3 มิติได้อยู่ในรูปแบบของ (148R, 52G, 200B)

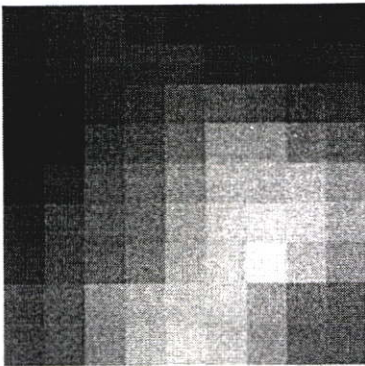
กำหนดจำนวนรอบในการเรียนรู้ $T=1000$ กำหนดรูปแบบของโหนดใกล้เคียงเป็นแบบสี่เหลี่ยม รัศมีของโหนดใกล้เคียง $\sigma(0)=5$ และอัตราเรียนรู้เริ่มต้น $\alpha(0)=0.2$ ในรูปที่ 3.7 ก. แสดงแผนภาพในตอนเริ่มต้นซึ่งจะทำการสุ่มค่าเวกเตอร์น้ำหนักในทันทีสุ่มเลือกตั้งแต่ 0-50 ทั้ง RGB รูปที่ 3.7 ข. รูปที่ 3.7 ค. รูปที่ 3.7 ง. แสดงแผนภาพ ณ จำนวนรอบที่ 100 250 และ 1000 รอบตามลำดับ



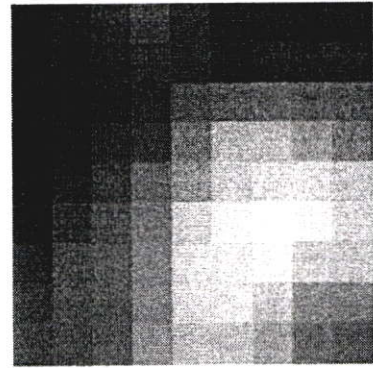
ก. แผนภาพเริ่มต้น



ข. แผนภาพ ณ จำนวน 100 รอบ



ค. แผนภาพ ณ จำนวน 250 รอบ



ง. แผนภาพ ณ จำนวน 1000 รอบ

รูปที่ 3.7 แสดงตัวอย่างแผนภาพ SOM ขนาด 9×9 ณ จำนวนรอบที่แตกต่างกัน

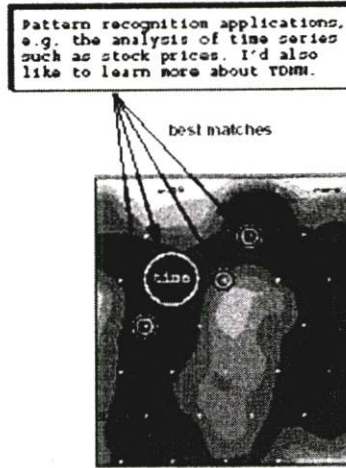
3.3 คุณสมบัติของ SOM

3.3.1 แผนภาพเรียงตัว

คุณสมบัติที่สำคัญที่ทำให้ SOM เป็นที่นิยมใช้งานกันอย่างแพร่หลายคือ คุณสมบัติในการวิเคราะห์ข้อมูลที่มีมิติสูง โดยผลลัพธ์ที่ได้จะถูกแสดงอยู่ในรูปแบบของแผนภาพ 2 มิติ เมื่อข้อมูลถูกกำหนดลงไปโหนดต่าง ๆ ของแผนภาพ ข้อมูลที่คล้ายกันจะถูกกำหนดให้โหนดที่อยู่ใกล้เคียงกัน ด้วยเหตุนี้แผนภาพที่ได้ออกมาจะมีลักษณะของการจัดเรียงตัวกันของข้อมูล ซึ่งทำให้ผู้ใช้สามารถที่จะเข้าใจลักษณะโครงสร้างของข้อมูลได้ นอกจากนี้การแสดงผลข้อมูลด้วยแผนภาพยังช่วยให้ผู้ใช้สามารถวิเคราะห์ข้อมูล มองเห็นความสัมพันธ์ของข้อมูล ซึ่งในบางครั้งไม่สามารถมองเห็นได้ด้วยการแสดงข้อมูลทั่วไปเช่น ตาราง หรือกราฟ

เข้ารหัสเป็นแผนภาพความหมาย (Semantic Map) จากนั้นสร้างเป็นแผนภาพของเอกสาร ดังรูปที่

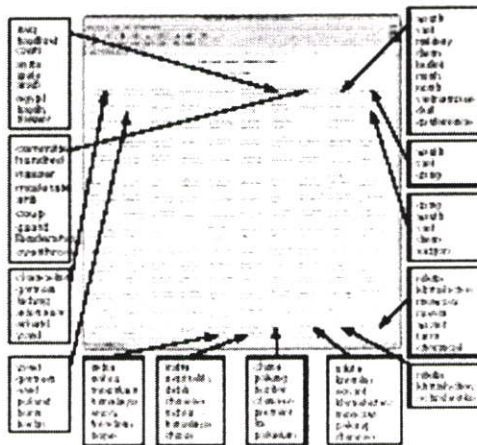
3.9



รูปที่ 3.9 แสดงแผนภาพ SOM ในงานวิจัย WEBSOM

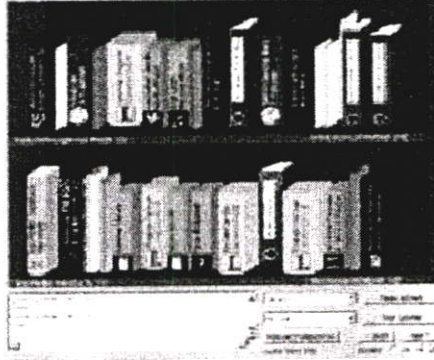
งานวิจัยนี้เป็นงานวิจัยแรก ๆ ที่นำเสนอลักษณะของ SOM ออกมาเป็นแผนภาพพร้อมสร้างอินเทอร์เน็ตในการค้นหา ถือเป็นงานนำเสนอการค้นคืนข้อมูลแบบใหม่

งานวิจัยของ Andreas Rauber, Dieter Merkl [12] ได้ประยุกต์ใช้งาน SOM เพื่อสร้างห้องสมุดอิเล็กทรอนิกส์(Digital Library) โดยนำเอาเอกสารมาสร้างเป็นอินเด็กซ์เทอมในรูปของ $tf \times idf$ [13,14] ซึ่งเป็นวิธีการที่ใช้ในวิทยานิพนธ์ฉบับนี้ ซึ่งจะกล่าวต่อไปในหัวข้อ 4.2 แผนภาพที่ใช้มีขนาด 10×15 ดังแสดงในรูปที่ 3.10 กระบวนการเรียนรู้จะเหมือนกับอัลกอริทึมที่ได้กล่าวในหัวข้อ 3.2



รูปที่ 3.10 แสดงการจัดกลุ่มเอกสาร โดยใช้แผนภาพ SOM ขนาด 10×15

ในงานวิจัยนี้ผู้วิจัยได้สร้างอินเตอร์เฟซที่ชื่อว่า LibViewer ขึ้นมาเพื่อจำลองภาพของชั้นวางหนังสือเสมือนขึ้นมามีลักษณะดังรูป 3.10 ในงานวิจัยนี้ได้แสดงให้เห็นถึงการนำแผนภาพ SOM ประยุกต์ใช้ในงานห้องสมุดเสมือนซึ่งเป็นอีกแนวคิดในการประยุกต์ใช้ SOM ในการจัดกลุ่มเอกสารเพื่ออำนวยความสะดวก

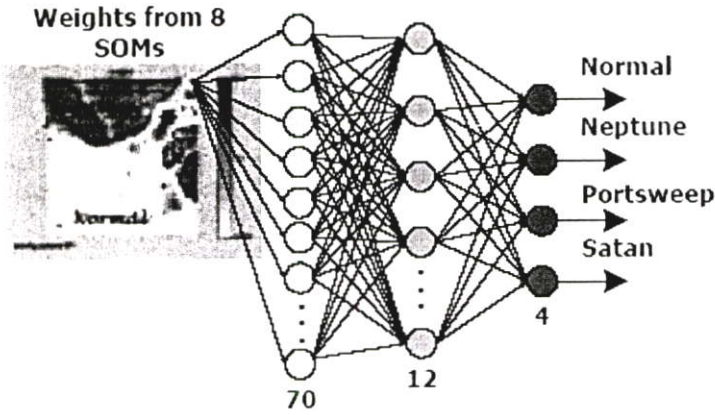


รูปที่ 3.11 แสดงชั้นวางหนังสือเสมือนใน LibViewer

ในปี 1990 Fox [15] ได้นำเสนอการประยุกต์ใช้ SOM ในการเรียนรู้ลักษณะพฤติกรรมการใช้งานปกติแล้วพิจารณาพฤติกรรมที่ผิดไปจากปกติว่าเป็นไวรัส

ในปี 1998 Cannady [16] ได้นำเสนอการใช้แผนภาพ SOM ประยุกต์ในการตรวจจับผู้บุกรุกแบบ Misuse detection

ในงานวิจัยของ Jirapummin C [5] ได้นำเอา SOM มาประยุกต์ใช้ร่วมกับ Resilient Propagation Neural Network (RPROP) สำหรับการตรวจจับผู้บุกรุก โดยใช้ SOM ในการจัดกลุ่ม (Clustering) พฤติกรรมผู้บุกรุกและแสดงบนแผนภาพ ส่วน RPROP ใช้สำหรับการแยกประเภทพฤติกรรมปกติกับพฤติกรรมผู้บุกรุก และทำการทดลองบนข้อมูลของ KDD Cup 1999 [x] ซึ่งเป็นข้อมูลที่จำลองพฤติกรรมผู้บุกรุก ซึ่งประกอบไปด้วย พฤติกรรมปกติ และพฤติกรรมการบุกรุก 3 ชนิด ได้แก่ การบุกรุกประเภท SYN flood คือ การบุกรุกชนิด neptune การบุกรุกประเภท port scan คือ การบุกรุกชนิด portsweep และการบุกรุกประเภท probe คือ การบุกรุกชนิด satan และแบ่งออกเป็นข้อมูลออกเป็น 8 กลุ่มย่อย ให้กับ SOM ทำการเรียนรู้ในแต่ละกลุ่มย่อยจนได้แผนภาพของ SOM แล้วจึงนำข้อมูลในแผนภาพส่งให้ RPROP 3 เลเยอร์ดังรูปที่ 3.12 และทำการแยกประเภทพฤติกรรมผู้บุกรุก ซึ่งได้ผลการทดลองออกมาดังตารางที่ 3.2



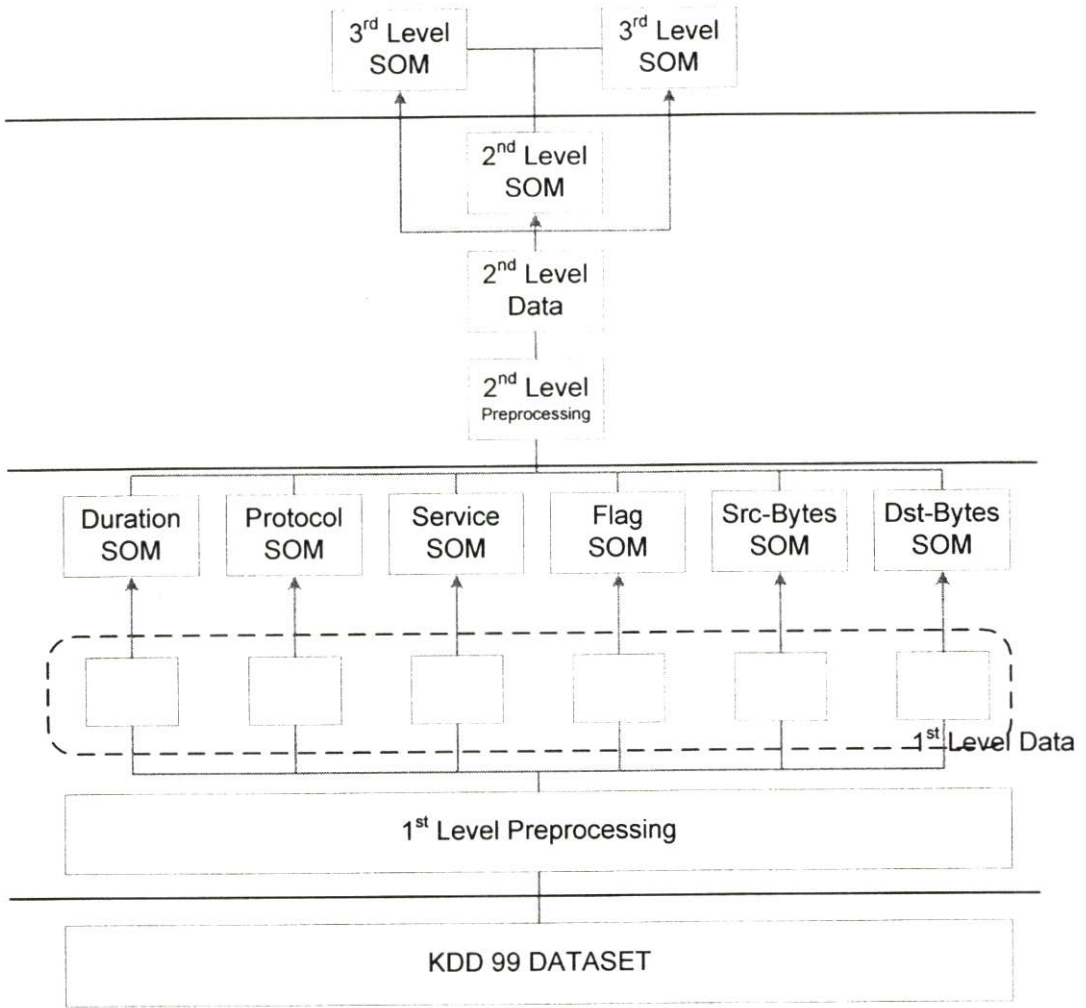
รูปที่ 3.12 แสดงโมเดลการทดลองของงานวิจัย

ผลของการตรวจจับผู้บุกรุกจากงานวิจัย [5] มีค่า Detection Rate และค่า False Alarm Rate ดีที่สุดคือ Neptune แต่โดยรวมระบบก็มีประสิทธิภาพดีดังตารางที่ 3.2

ตารางที่ 3.2 แสดงผลการทดลองของงานวิจัย

Attacks	Detection Rate	False Alarm Rate
Neptune	99.7181	0.0591
Port Sweep	97.9123	4.1917
Satan	90.2811	4.4988

ในปี 2003 Heywood [6] ได้นำเสนอ SOM ในการตรวจจับผู้บุกรุกโดยใช้ Toolbox จาก MATLAB ในการสร้างแผนภาพ SOM โดยทำการทดลองบนข้อมูล KDD Cup 1999 สำหรับการเรียนรู้และทดสอบโดยออกแบบ SOM เป็น 3 ลำดับชั้น โดยมีวิธีการดังนี้ คือ เริ่มโดยกำหนดเลือกข้อมูลจาก KDD Cup 1999 มาเพียง 6 คุณลักษณะจาก 41 คุณลักษณะ ได้แก่ duration, protocol, service, flag, destination และ source ซึ่งในลำดับชั้นที่ 1 จะทำการสร้างแผนภาพ SOM ออกมาเป็น 6 แผนภาพเพื่อสำหรับเรียนรู้ข้อมูลคุณลักษณะทั้ง 6 คุณลักษณะ ลำดับชั้นที่ 2 จะสร้างแผนภาพ SOM จำนวน 1 แผนภาพรวมข้อมูลจากแผนภาพทั้ง 6 คุณลักษณะในลำดับชั้นที่ 1 โดยใช้ฟังก์ชันในการจัดเรียงกลุ่มข้อมูล (Potential function clustering) ของจำนวนข้อมูลอินพุต และสำหรับการหาโหนดใกล้เคียงของ SOM ใช้ฟังก์ชันเกาส์ โดยมีโมเดลดังรูปที่ 3.13



รูปที่ 3.13 แสดงโมเดลการทดลองของงานวิจัย [6]

ผลที่ได้จากการทดลองของ Heywood คือ ได้ 89 % Detection rate และ False positive rate ที่ 4.6%

3.5 สรุป

ในบทนี้ได้นำเสนอพื้นฐานของการตรวจจับผู้บุกรุก รูปแบบการเรียนรู้พฤติกรรมผู้บุกรุก โครงสร้างของ SOM อัลกอริธึมการเรียนรู้ และคุณสมบัติของ SOM พร้อมทั้งแสดงตัวอย่างการใช้งาน เพื่อให้เห็นประโยชน์ของการใช้งาน SOM รวมถึงงานวิจัยที่เกี่ยวข้อง ในบทถัดไปจะนำเสนอปัญหาที่พบในการใช้งาน SOM ในการตรวจจับผู้บุกรุก พร้อมทั้งนำเสนอโมเดลใหม่

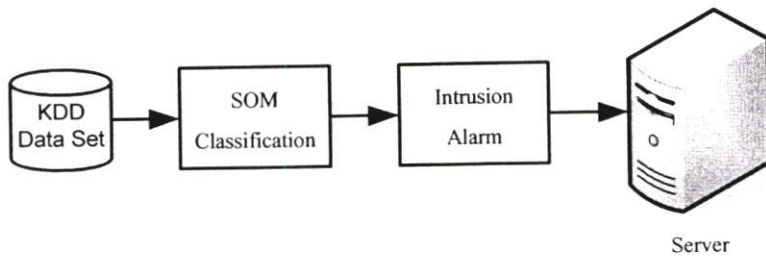
บทที่ 4

การออกแบบมัลติเลเยอร์ SOM

ในบทนี้จะกล่าวถึงหลักการและวิธีการดำเนินงานวิจัย ของงานวิจัยที่ผู้วิจัยนำเสนอ ซึ่งในส่วนนี้จะนำเสนอถึงวิธีการในการดำเนินงานวิจัย การเตรียมข้อมูลที่ใช้ในการทดลอง แสดงโมเดลการทดลอง และการวัดประสิทธิภาพของระบบ เป็นต้น

4.1 โมเดลการทำงานของระบบ

ในหัวข้อนี้แนะนำขั้นตอนการทำงานของโมเดลระบบที่ผู้วิจัยได้นำเสนอ ซึ่งเป็นการประยุกต์ใช้ SOM แบบหลายลำดับชั้น หรือเรียกว่า “Multi-layer Self-Organizing Map” ในการจำแนกประเภทการบุกรุกของข้อมูล ดังรูปที่ 4.1



รูปที่ 4.1 แสดงขั้นตอนการทำงานของระบบ

โดยในส่วนแรก คือ ข้อมูลที่ใช้ในงานวิจัยนี้เป็นข้อมูลที่ได้จาก KDD Cup 1999 [17] ซึ่งเป็นข้อมูลที่จำลองพฤติกรรมกรรมการบุกรุกที่นิยมใช้ในงานวิจัยที่เกี่ยวกับการตรวจจับผู้บุกรุกจะแสดงรายละเอียดในหัวข้อถัดไป ต่อมาในส่วนที่สองจะเป็นส่วนของการประยุกต์ใช้ SOM มาช่วยในการจำแนกประเภทการบุกรุกเครือข่ายคอมพิวเตอร์ ซึ่งในงานวิจัยนี้ได้แนะนำการประยุกต์ใช้ SOM แบบหลายลำดับชั้นมาช่วยในการตรวจจับผู้บุกรุก ส่วนที่สามจะเป็นส่วนของการแจ้งเตือนไปสู่เครื่องคอมพิวเตอร์ และในส่วนสุดท้ายจะเป็นส่วนที่ผู้ดูแลระบบจะได้รับข้อมูลจากการจำแนกประเภทพฤติกรรมกรรมการบุกรุกจาก SOM เพื่อเก็บไว้เป็นข้อมูลในเครื่องคอมพิวเตอร์แม่ข่าย

4.2 ข้อมูล KDD Cup 1999

ข้อมูล KDD Cup 1999 [17] เป็นข้อมูลที่ประยุกต์มาจากข้อมูลพฤติกรรมกรรการบุกรุกของ DARPA 98 [17] ซึ่งมีข้อมูลประกอบไปด้วยทั้งพฤติกรรมกรรการบุกรุกและพฤติกรรมปกติ โดยข้อมูลนี้ได้จากการวิเคราะห์ระบบเครือข่ายของมหาวิทยาลัยโคลัมเบีย (Columbia University) [17]

4.2.1 ลักษณะข้อมูลของ KDD Cup 1999

เป็นข้อมูลของการติดต่อสื่อสารบนโปรโตคอล TCP ในระบบเครือข่ายคอมพิวเตอร์ ซึ่งประกอบไปด้วยพฤติกรรมปกติ และพฤติกรรมกรรการบุกรุก โดยพฤติกรรมกรรการบุกรุกได้แบ่งออกเป็น 4 กลุ่มใหญ่ๆ ดังต่อไปนี้ [17]

1. Denial of Service เป็นลักษณะของผู้บุกรุกพยายามโจมตีระบบคอมพิวเตอร์ไม่ให้สามารถให้บริการต่างๆ ได้ เช่น smurf
2. Remote to Local เป็นลักษณะของผู้บุกรุกที่ไม่ได้เป็นยูสเซอร์ในระบบแต่พยายามเจาะเข้าไปในระบบ เช่น guess password
3. User to Root เป็นลักษณะของผู้บุกรุกที่พยายามเข้าสู่ระบบโดยการใช้สิทธิ์ของซูเปอร์ยูสเซอร์ เช่น buffer overflow
4. Probe เป็นลักษณะของผู้บุกรุกที่พยายามตรวจสอบหาจุดอ่อนของระบบ เช่น portsweep

และในข้อมูลของ KDD Cup 1999 จะประกอบไปด้วยคุณลักษณะจำนวน 41 คุณลักษณะที่เป็นคุณลักษณะที่ได้จากการเชื่อมต่อระบบเครือข่าย โดยสามารถแยกคุณลักษณะออกเป็นกลุ่มย่อย ประกอบไปด้วย คุณลักษณะพื้นฐาน (Basic features) มีจำนวนทั้งสิ้น 9 คุณลักษณะ และคุณลักษณะเพิ่มเติมอีก 3 กลุ่ม โดยมีจำนวนทั้งสิ้น 32 คุณลักษณะ ดังต่อไปนี้ [17]

1. Basic features เป็นคุณลักษณะพื้นฐานที่ได้จากแพคเกจข้อมูลที่สื่อสารในเครือข่าย ประกอบไปด้วย

ตารางที่ 4.1 คุณลักษณะพื้นฐาน

Feature name	Description
Duration	Length (number of seconds) of the connection
Protocol_type	Type of the protocol, e.g., tcp, udp, etc.
Service	Network service on the destination, e.g., http, telnet, etc.
Flag	Normal or error status of the connection
Src_bytes	Number of data bytes from source to destination
Dst_bytes	Number of data bytes from destination to source

ตารางที่ 4.1 (ต่อ)

Land	“1” if connection is from/to the same host/port; “0” otherwise
Wrong_fragment	Number of “wrong” fragments
Urgent	Number of urgent packet

2. Content features เป็นคุณลักษณะที่เก็บรวบรวมข้อมูลที่แสดงให้เห็นถึงพฤติกรรมน่าสงสัย เช่น ความผิดพลาดในการล็อกอิน หรือการใช้คำสั่ง “su” เป็นต้น ประกอบไปด้วย

ตารางที่ 4.2 Content features

Feature name	Description
Hot	Number of “hot” indicators
Num_failed_logins	Number of failed login attempts
Logged_in	“1” if successfully logged in; “0” otherwise
Num_compromised	Number of “compromised” conditions
Root_shell	“1” if root shell is obtained; “0” otherwise
Su_attempted	“1” if “su root” command attempted; “0” otherwise
Num_root	Number of “root” accesses
Num_file_creations	Number of file creation operations
Num_shells	Number of shell prompts
Num_access_files	Number of operations on access control files
Num_outbound_cmds	Number of outbound commands in an ftp session
Is_host_login	“1” if the login belongs to the “hot” list; “0” otherwise
Is_guest_login	“1” if the login is a “guest” login; “0” otherwise

3. Traffic features เป็นคุณลักษณะที่เก็บรวบรวมข้อมูลที่แสดงลักษณะของการสื่อสาร เช่น จำนวนครั้งในการเชื่อมต่อเข้าสู่ระบบเมื่อผ่านไป 2 วินาที เป็นต้น ประกอบไปด้วย

ตารางที่ 4.3 Traffic feature

Feature name	Description
Count	Number of connections to the same host as the current connection in the past two seconds.

ตารางที่ 4.3 (ต่อ)

	Note: The following features refer to these same-host connection.
Srv_count	number of connections having the same service as the connection.
Serror_rate	S0 error rate
Srv_serror_rate	S0 error rate for the same service as the current one.
Rerror_rate	RST error rate
Srv_rerror_rate	RST error rate for the same service as the current one.
Same_srv_rate	Percentage of connections that use the same service as the current one.
Diff_srv_rate	Percentage of different service.
Srv_diff_host_rate	Percentage of different host used by the current service.

4. Host based features เป็นคุณลักษณะที่เก็บรวบรวมข้อมูลที่แสดงลักษณะของการสื่อสารไปยังเครื่องปลายทางเครื่องเดิมตลอดเวลา เช่น จำนวนครั้งในการเชื่อมต่อไปยังเครื่องปลายทางเครื่องเดิม เป็นต้น ประกอบไปด้วย

ตารางที่ 4.4 Host based feature

Feature name	Description
Dst_host_count	Count of connections having the same destination host.
Dst_host_srv_count	Count of connections having the same destination host and using the same service.
Dst_host_same_srv_rate	Percentage of connections having the same destination host and using the same service.
Dst_host_diff_srv_rate	Percentage of different services on the current host.
Dst_host_same_src_port_rate	Percentage of connections to the current host having the same src port.
Dst_host_srv_diff_host_rate	Percentage of connections to the same service coming from different hosts.
Dst_host_serror_rate	Percentage of connections to the current host that have an S0 error.
Dst_host_srv_serror_rate	Percentage of connections to the current host and specified service that have an S0 error.

ตารางที่ 4.4 (ต่อ)

Dst_host_error_rate	Percentage of connections to the current host that have an RST error.
Dst_host_srv_error_rate	Percentage of connections to the current host and specified service that have an RST error.

4.2.2 การเตรียมข้อมูลอินพุตเวกเตอร์

ในการดำเนินงานวิจัยนี้ผู้ดำเนินงานวิจัยได้เตรียมข้อมูลอินพุตสำหรับระบบ เพื่อใช้ในการทดลอง โดยข้อมูลที่ใช้ในการทดลองกำหนดบนพื้นฐานของประเภทการบุกรุกแบบ Denial of service กับ Probing และพฤติกรรมปกติ เนื่องจากเป็นกลุ่มข้อมูลที่มีจำนวนมากที่สุด มีทั้งสิ้น 492,843 เรคคอร์ด และสำหรับคุณลักษณะทั้ง 41 คุณลักษณะ จะมีคุณลักษณะบางคุณลักษณะนั้นเป็นคุณลักษณะแบบสัญลักษณ์ (Symbolic) ซึ่งต้องทำการจัดให้อยู่ในลักษณะของตัวเลขก่อน เพื่อให้เป็นไปตามข้อกำหนดพื้นฐานของข้อมูลอินพุตในการนำไปใช้ใน SOM โดยมีรายละเอียดดังต่อไปนี้

ตารางที่ 4.5 แสดงการกำหนดค่าตัวเลขแทน Protocol feature

Value	Assigned
tcp	0
udp	1
icmp	2

ตารางที่ 4.6 แสดงการกำหนดค่าตัวเลขแทน Service feature

Value	Assigned	Value	Assigned
http	0	other	10
smtp	1	private	11
finger	2	pop_3	12
domain_u	3	ftp_data	13
auth	4	rje	14
telnet	5	time	15
ftp	6	mtp	16
eco_i	7	link	17
ntp_u	8	remote_job	18
ecr_i	9	gopher	19

ตารางที่ 4.6 แสดงการกำหนดค่าตัวเลขแทน Service feature (ต่อ)

Value	Assigned	Value	Assigned
ssh	20	iso_tsap	44
name	21	hostname	45
whois	22	csnet_ns	46
domain	23	pop_2	47
login	24	sunrpc	48
imap4	25	uucp_path	49
daytime	26	netbios_ns	50
ctf	27	netbios_ssn	51
nntp	28	netbios_dgm	52
shell	29	sql_net	53
IRC	30	vmnet	54
nnspp	31	bgp	55
http_443	32	Z39_50	56
exec	33	tim_i	64
printer	34	red_j	65
efs	35	ldap	57
courier	36	netstat	58
uucp	37	urh_i	59
klogin	38	X11	60
kshell	39	urp_i	61
echo	40	pm_dump	62
discard	41	tftp_u	63
sysstat	42	tim_i	64
supdup	43	red_j	65

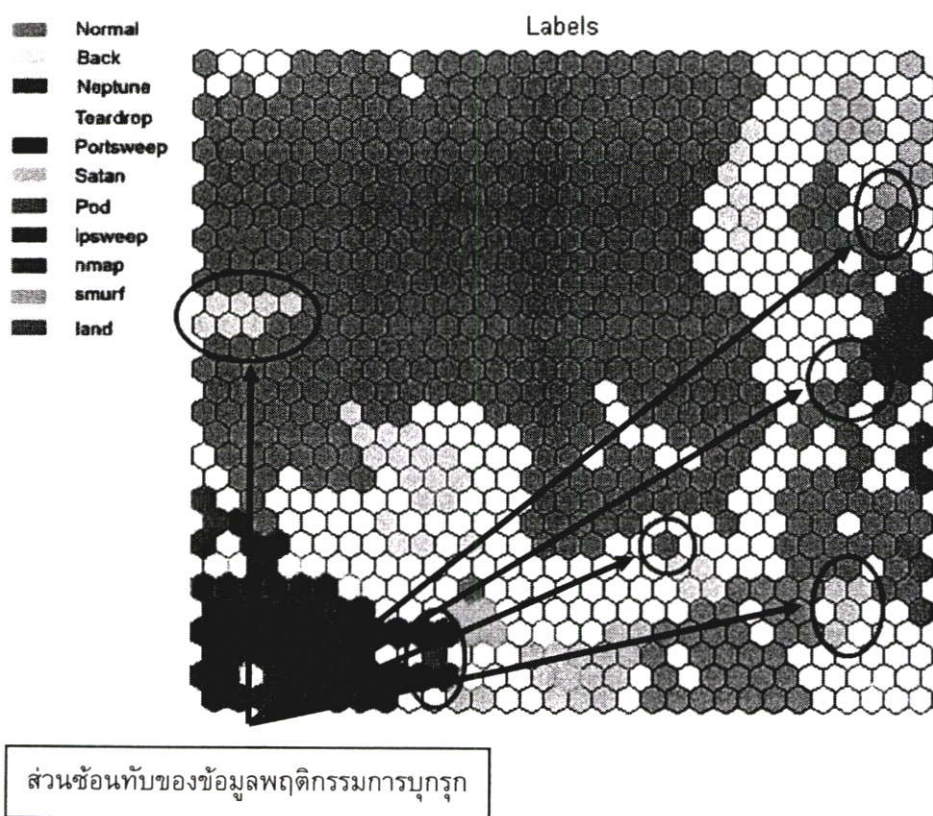
ตารางที่ 4.7 แสดงการกำหนดค่าตัวเลขแทน Flag feature

Value	Assigned	Value	Assigned	Value	Assigned
SF	0	S0	4	RSTOS0	8
S1	1	S3	5	OTH	9
REJ	2	RSTO	6	SH	10
S2	3	RSTR	7		

จากนั้นจึงแบ่งข้อมูลออกเป็น 2 กลุ่มเพื่อใช้ในการสอนระบบจำนวน 394,274 เรคคอร์ด และใช้ในการทดสอบระบบจำนวน 98,569 เรคคอร์ด โดยมีพฤติกรรมการบุกรุกจำนวน 11 ชนิด

4.3 การสร้างแผนภาพ SOM ลำดับชั้นที่หนึ่ง

การทดลองในงานวิจัยนี้ผู้วิจัยได้ทำการทดลองเพื่อหาแผนภาพ SOM ในลำดับชั้นที่ 1 โดยพยายามหาแผนภาพ SOM ที่สามารถแสดงประเภทการบุกรุก ที่มีการทับซ้อนกันน้อยที่สุด และจากผลการทดสอบในการหาแผนภาพ SOM ดังกล่าว เราจึงได้แผนภาพ SOM ที่แสดงประเภทการบุกรุกที่มีการทับซ้อนกันน้อยที่สุด คือขนาด 30x30 หรือขนาด 900 โหนดเพื่อใช้ในแผนภาพ SOM ชั้นแรก ดังรูปที่ 4.2



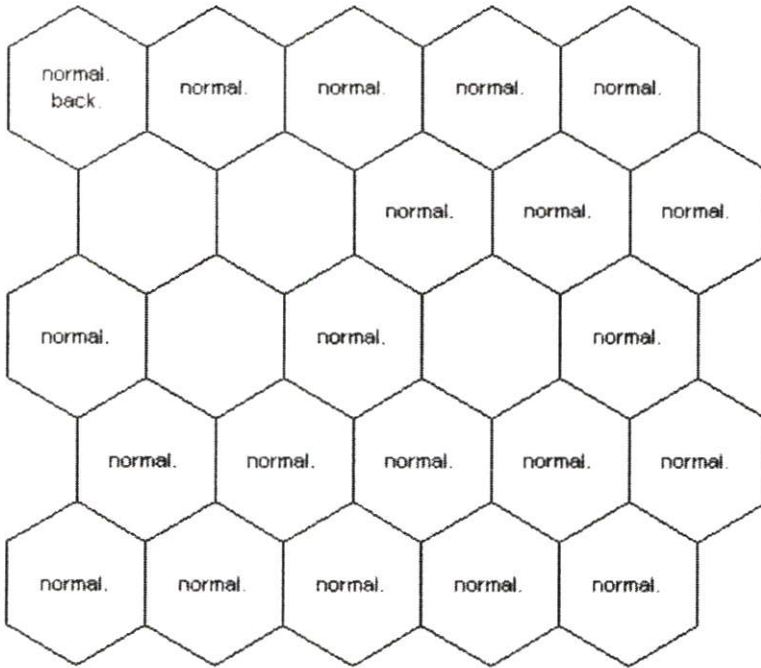
รูปที่ 4.2 แสดงแผนภาพ SOM ขนาด 900 โหนด หรือขนาด 30 x 30

4.4 แผนภาพ SOM ลำดับชั้นต่อไป

ในงานวิจัยนี้ผู้วิจัยจะพิจารณาปรับปรุงข้อจำกัดเรื่องของการทับซ้อนกันของพฤติกรรมผู้บุกรุกโดยนำเสนอวิธีการใหม่ ซึ่งอาศัยแนวคิดดังนี้ คือ เมื่อพิจารณาความหมายของการทับซ้อนกันของข้อมูล การทับซ้อนกันของข้อมูลนั้นเกิดจากการที่มีข้อมูลที่มีเนื้อหาใกล้เคียง

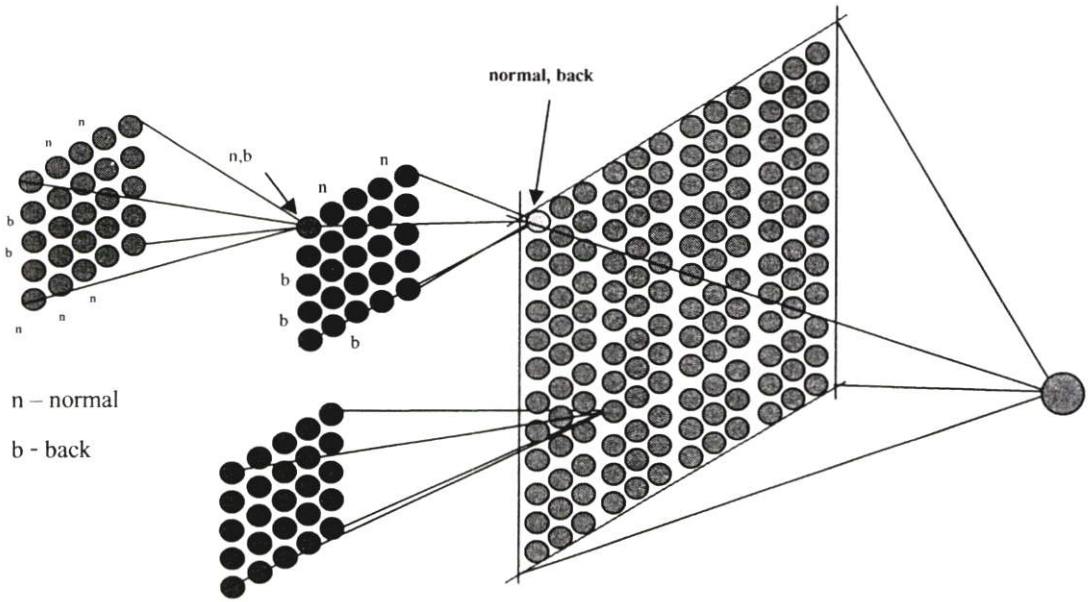
กัน ตกอยู่ในโหนดขณะเดียวกันซึ่งเป็นโหนดที่เป็นตัวแทนกลุ่มข้อมูลที่มีเนื้อหาสอดคล้องกับข้อมูลอินพุต ดังแสดงในรูปที่ 4.3 และเมื่อเป็นดังนี้เราจึงเสนอให้มีการสร้างแผนภาพชั้นใหม่อีกชั้นหนึ่ง เพื่อกระจายข้อมูลที่ทับซ้อนกัน โดยกำหนดค่าคงที่ค่าหนึ่งเพื่อใช้เป็นข้อกำหนดของการสร้างแผนภาพชั้นใหม่ต่อไป เป้าหมายจะทำการสร้างแผนภาพไปเรื่อยๆ จนกว่าค่าอัตราส่วนของประเภทการบุกรุกในโหนดที่มีข้อมูลทับซ้อนกันมีค่ามากกว่าหรือเท่ากับค่าคงที่ที่กำหนดจึงหยุดสร้างแผนภาพชั้นใหม่ จากหลักการที่ได้กล่าวมา เราได้นำมาสร้างเป็นแผนภาพดังแสดงในรูปที่

4.4



Trian Data 25 Node

รูปที่ 4.3 แสดงตัวอย่างข้อมูลที่ทับซ้อนในแผนภาพ SOM



รูปที่ 4.4 แสดงการแบ่งชั้นของแผนภาพ SOM หลายลำดับชั้น

จากรูปที่ 4.4 แสดงการแบ่งชั้นของแผนภาพ SOM หลายลำดับชั้น โดยที่ชั้นที่ 1 ซึ่งเป็นชั้นเอาต์พุต ของการเรียนรู้ของ SOM จะถูกดำเนินการตรวจสอบข้อมูลที่ทับซ้อนกัน และจะถูกนำไปเข้าสู่กระบวนการตรวจสอบหาค่า อัตราส่วนของชนิดการบุกรุก (Class Ratio) ดังสมการที่ (4.1)

$$CR = \left(\frac{N_{Select}}{N_{Tn}} \right) \quad (4.1)$$

เมื่อ	CR	แทนค่าอัตราส่วนชนิดของการบุกรุก
	N_{Select}	แทนค่าจำนวนประเภทการบุกรุกที่เลือกในโหนด
	N_{Tn}	แทนค่าจำนวนประเภทการบุกรุกทั้งหมดใน โหนด

และนำค่า CR ที่ได้มาเปรียบเทียบกันในแต่ละโหนดเพื่อหาค่าสูงสุด (CR_{Max}) และนำค่าอัตราส่วนสูงสุด ไปเปรียบเทียบค่าคงที่ค่าหนึ่ง เพื่อใช้ในการพิจารณาว่าควรเพิ่มชั้นของแผนภาพ SOM หรือไม่ ซึ่งสามารถแสดงอัลกอริทึมได้ในหัวข้อถัดไป

4.5 อัลกอริทึมในการเพิ่มชั้นของแผนภาพ SOM

อัลกอริทึมในการเพิ่มชั้นของแผนภาพ SOM สามารถแสดงได้ดังนี้

1. หาโหนดในแผนภาพชั้นที่หนึ่ง ที่มีข้อมูลที่ทับซ้อนกัน

2. กำหนดค่าอัตราส่วน CR ดังสมการที่ (3.1) ในแต่ละประเภทของการบุกรุกและเปรียบเทียบหาค่าสูงสุด(CR_{Max})
3. พิจารณาเพิ่มลำดับชั้น ดังสมการที่ (3.2)
 - 3.1. ถ้าค่าอัตราส่วนสูงสุด(CR_{Max}) มากกว่าหรือเท่ากับ ค่าที่กำหนด(β) จะไม่เพิ่มลำดับชั้น
 - 3.2. ถ้าค่าอัตราส่วนสูงสุด(CR_{Max}) น้อยกว่า ค่าที่กำหนด(β) จะเพิ่มลำดับชั้น

$$AddLayer(CR_{Max}) = \begin{cases} true; & \text{if } CR_{Max} < \beta \\ false; & \text{otherwise} \end{cases} \quad (4.2)$$

4. ทำซ้ำในข้อที่ 1 ถึง 3 ตามลำดับจนครบทุกโหนดที่มีข้อมูลที่ทับซ้อน

4.6 การวัดค่า Detection Rate และ False Positive

ในโลกของความเป็นจริงการตรวจจับการบุกรุกของระบบที่ IDS รับรู้อาจมีความเหลื่อมล้ำกับระบบที่ IDS ควรตรวจจับได้จริง ทำให้เกิดความผิดพลาดในการตรวจจับการบุกรุกได้ ซึ่งเราสามารถแบ่งความผิดพลาด ในการตรวจจับการบุกรุกได้เป็นสองลักษณะคือ false positive และ false negative ซึ่งความผิดพลาดทั้งสองแบบมีรายละเอียด คือ

1. false Positive คือ ความผิดพลาดอันเนื่องมาจากการเกิดเหตุการณ์ซึ่งเป็นเหตุการณ์ปกติในระบบ แต่ IDS คิดว่าเกิดเหตุการณ์ผิดปกติเกิดขึ้น ผลลัพธ์คือ IDS จึงแจ้งเตือนผู้ดูแลระบบว่าเกิดเหตุการณ์ผิดปกติ
2. false Negative คือ ความผิดพลาดอันเนื่องมาจากการเกิดเหตุการณ์ซึ่งเป็นเหตุการณ์ผิดปกติในระบบ แต่ IDS คิดว่าเป็นเหตุการณ์ปกติ จึงไม่ได้แจ้งเตือนผู้ดูแลระบบว่าเกิดเหตุการณ์ผิดปกติ

โดยการออกแบบระบบ IDS นั้นจะพยายามออกแบบ และใช้วิธีต่างๆ มากมายที่ทำให้ false positive และ false negative มีน้อยที่สุด แต่ในงานวิจัยนี้เป็นงานวิจัยที่ดำเนินการ โดยเปรียบเทียบกับค่าของพฤติกรรมปกติ เราจึงให้ความสนใจในค่าของ false Positive เพียงปัจจัยเดียวดังสมการที่ (4.3)

$$\% \text{ false_Positive} = \left(\frac{N_{false}}{N_{normal}} \right) \times 100 \quad (4.3)$$

เมื่อ

- N_{false} แทนค่าจำนวนทั้งหมดของ false positive ที่พบในการทดสอบ
- N_{normal} แทนค่าจำนวนพฤติกรรมปกติ (normal) ทั้งหมดในข้อมูลทดสอบ

และมีการวัดค่าความถูกต้องในการตรวจจับ หรือที่เรียกว่า “detection rate” ดังสมการที่ (4.4)

$$\%Detected = \left(\frac{N_{attack} - N_{missed}}{N_{attack}} \right) \times 100 \quad (4.4)$$

เมื่อ

N_{attack} แทนค่าจำนวนของประเภทการบุกรุกทั้งหมดในข้อมูลทดสอบ

N_{missed} แทนค่าจำนวนของการตรวจจับผิดพลาด

4.7 สรุป

ในบทนี้ผู้วิจัยนำเสนอวิธีการออกแบบมัลติเลเยอร์ SOM เพื่อใช้ในการตรวจจับการบุกรุกระบบเครือข่ายคอมพิวเตอร์ ซึ่งการออกแบบมัลติเลเยอร์ SOM นี้เป็นแนวความคิดที่ได้จากการพยายามนำเสนอรูปแบบพฤติกรรมการบุกรุกระบบเครือข่ายคอมพิวเตอร์ให้สามารถจำแนกและแยกประเภทพฤติกรรมการบุกรุกออกมาได้อย่างชัดเจน เพื่อเป็นประโยชน์ต่อผู้ดูแลระบบเครือข่ายคอมพิวเตอร์ โดยลักษณะเด่นของ SOM นี้อยู่ตรงการนำเสนอแผนภาพการจัดกลุ่มของข้อมูล ทางผู้วิจัยจึงได้นำลักษณะเด่นของ SOM นี้มาช่วยในการนำเสนอแผนภาพการจัดกลุ่มของข้อมูลพฤติกรรมการบุกรุก

เมื่อผู้วิจัยได้ทำการประยุกต์ SOM มาช่วยในการตรวจจับการบุกรุกระบบเครือข่ายในเบื้องต้นผู้วิจัยได้พบว่า SOM สามารถนำเสนอแผนภาพที่จัดกลุ่มข้อมูลพฤติกรรมการบุกรุก และแยกประเภทพฤติกรรมออกเป็นกลุ่มๆ ได้อย่างมีประสิทธิภาพในระดับหนึ่ง เพราะว่าผู้วิจัยได้ทำการทดสอบโดยการนำ SOM หลายขนาดมาทำการทดสอบกับกลุ่มข้อมูลพฤติกรรมการบุกรุก ผลปรากฏว่าแผนภาพของ SOM นั้นสามารถนำเสนอและแยกประเภทพฤติกรรมการบุกรุกออกเป็นกลุ่มๆ ได้แตกต่างกันขึ้นอยู่กับขนาดของแผนภาพ SOM โดยขนาดที่เหมาะสมที่สุดผู้วิจัยเลือกจากการจัดกลุ่มข้อมูลพฤติกรรมและแยกประเภทการบุกรุกออกเป็นกลุ่มๆ ให้เห็นว่ามียังมีจำนวนกลุ่มเท่ากับจำนวนกลุ่มของพฤติกรรมการบุกรุกตัวอย่างที่นำมาทดสอบจำนวนทั้งสิ้น 11 กลุ่ม คือ ขนาด 30x30 หรือ 900 โหนด แต่ปัญหาในการทดสอบแผนภาพ SOM เบื้องต้นนี้ คือ การทับซ้อนกันของข้อมูลพฤติกรรมการบุกรุกตัวอย่างที่นำมาทดสอบใน โหนดเดียวกันของแผนภาพ

ดังนั้นผู้วิจัยจึงได้นำเสนอแนวความคิดในการสร้างแผนภาพ SOM ขึ้นซ้อนกับแผนภาพ SOM เดิม โดยคาดหวังว่าจะช่วยลดปัญหาการทับซ้อนของการจัดกลุ่มข้อมูลของแผนภาพ เพราะข้อมูลได้ถูกกระจายจัดกลุ่มใหม่ในแผนภาพใหม่ที่สร้างขึ้น และเมื่อผู้วิจัยได้ทำการทดสอบในบางโหนดของแผนภาพ ผลปรากฏว่าแผนภาพที่สร้างขึ้นใหม่ได้ช่วยแก้ปัญหาการทับซ้อนกันของข้อมูลพฤติกรรมได้อย่างมีประสิทธิภาพ แต่ปัญหาต่อมาในการสร้างแผนภาพใหม่ขึ้นมาซ้อนแผนภาพเดิมนั้น ถ้าหากว่าทำการกำหนดขนาดที่ใหญ่ คือ มีจำนวนโหนดมาก ก็จะมีการเสียเวลาในการคำนวณหาโหนดผู้ชนะมากขึ้น ผู้วิจัยจึงนำเสนอแผนภาพที่สร้างใหม่ให้มีขนาดเล็กที่สุด

เพื่อลดเวลาในการคำนวณของแผนภาพ ผลปรากฏว่าขนาดที่เหมาะสม คือ ขนาด 5x5 หรือจำนวน 25 โหนด โดยผู้วิจัยพิจารณาจากจำนวนของโหนดที่ไม่ได้เป็นตัวแทนของกลุ่มข้อมูลที่เรียกว่า “โหนดตาย” ซึ่งหากมีจำนวนโหนดตายมาก และเมื่อนับจำนวนโหนดที่เป็นตัวแทนของกลุ่มข้อมูลที่เรียกว่า “โหนดเป็น” ทั้งหมดของทุกแผนภาพรวมกันแล้วมีจำนวนมากที่สุด ก็จะไม่เหมาะสม และในขั้นตอนนี้ผู้วิจัยได้นำเสนอการกำหนดค่า β สำหรับใช้ในการพิจารณาว่า โหนดที่มีข้อมูลทับซ้อนอยู่นั้นสมควรสร้างแผนภาพ SOM ในชั้นใหม่หรือไม่ เพื่อลดปัญหาเรื่องเวลาในการคำนวณของแผนภาพ SOM เพราะจะสร้างแผนภาพ SOM ใหม่เฉพาะโหนดที่มีค่าน้อยกว่าค่าอัตราส่วนที่กำหนดเท่านั้น ไม่ใช่สร้างแผนภาพ SOM ใหม่สำหรับทุกโหนดของแผนภาพ SOM ในลำดับชั้นแรก

สุดท้ายผู้วิจัยได้นำแผนภาพทั้งหมดมาวัดประสิทธิภาพในการตรวจจับการบุกรุก เพื่อแสดงให้เห็นว่าแนวความคิดดังกล่าวมานี้เป็นแนวความคิดที่สามารถนำไปประยุกต์ใช้ได้โดยมีประสิทธิภาพต่อไป

บทที่ 5

ผลการทดลอง

ในบทนี้จะกล่าวถึงการนำโมเดลจากบทที่ 4 มาทดลองเพื่อทดสอบประสิทธิภาพของโมเดล SOM แบบลำดับชั้นเดียว เปรียบเทียบกับโมเดล SOM แบบหลายลำดับชั้น ในการตรวจจับผู้บุกรุก โดยการทดลองประกอบไปด้วย 2 การทดลอง การทดลองที่ 1 เป็นการทดสอบสร้างแผนภาพ SOM แบบลำดับชั้นเดียวในการตรวจจับผู้บุกรุก เพื่อหาแผนภาพ SOM ที่เหมาะสม การทดลองที่ 2 เป็นการทดสอบสร้างแผนภาพ SOM แบบหลายลำดับชั้นโดยใช้ค่าคงที่ β เท่ากับ 0.8, 0.9 และ 1.0 เพื่อแก้ปัญหาในการจำแนกประเภทข้อมูลการบุกรุกที่ทับซ้อนกันของระบบ SOM แบบลำดับชั้นเดียว

5.1 การทดลองที่ 1 การตรวจจับการบุกรุกโดยใช้ SOM แบบลำดับชั้นเดียว

5.1.1 จุดประสงค์การทดลอง

ในการทดลองนี้ผู้วิจัยมีวัตถุประสงค์ต้องการทดสอบการประยุกต์ใช้แผนภาพ SOM กับชุดข้อมูลของ KDD Cup 1999 เพื่อหาแผนภาพ SOM ที่เหมาะสมสำหรับการจัดกลุ่มพฤติกรรมการบุกรุกระบบเครือข่ายคอมพิวเตอร์

5.1.2 ขั้นตอนการทดลอง

ผู้วิจัยได้กำหนดสถานะในการทดลองดังตารางที่ 5.1 และกำหนดข้อมูลที่ใช้เป็นข้อมูลที่ได้จาก KDD Cup 1999

ตารางที่ 5.1 พารามิเตอร์ในการสอนระบบเริ่มต้น

Parameter	Rough Training	Fine Tuning
Initial α	0.5	0.05
α decay scheme	inverse_t	
Epoch Limit	4,000	
Neighborhood Parameter		
Initial Size	2	1
Function	Gussian	
Relation	Hexagonal	

ในการทดลองนี้ผู้วิจัยได้จัดเตรียมข้อมูลสำหรับทำการสอนระบบและทดสอบระบบ ดังตารางที่ 5.2

ตารางที่ 5.2 ข้อมูลที่ใช้ในการสอนและทดสอบระบบจำนวนทั้งสิ้น 492,843 เรคคอร์ด

Attack Type	Train	Test
normal.	77,822	19,456
neptune.	85,761	21,440
smurf.	224,632	56,158
back.	1,762	441
satan.	1,271	318
ipsweep.	998	249
portsweep.	832	208
teardrop.	783	196
pod.	211	53
land.	17	4
nmap.	185	46
Total	394,274	98,569

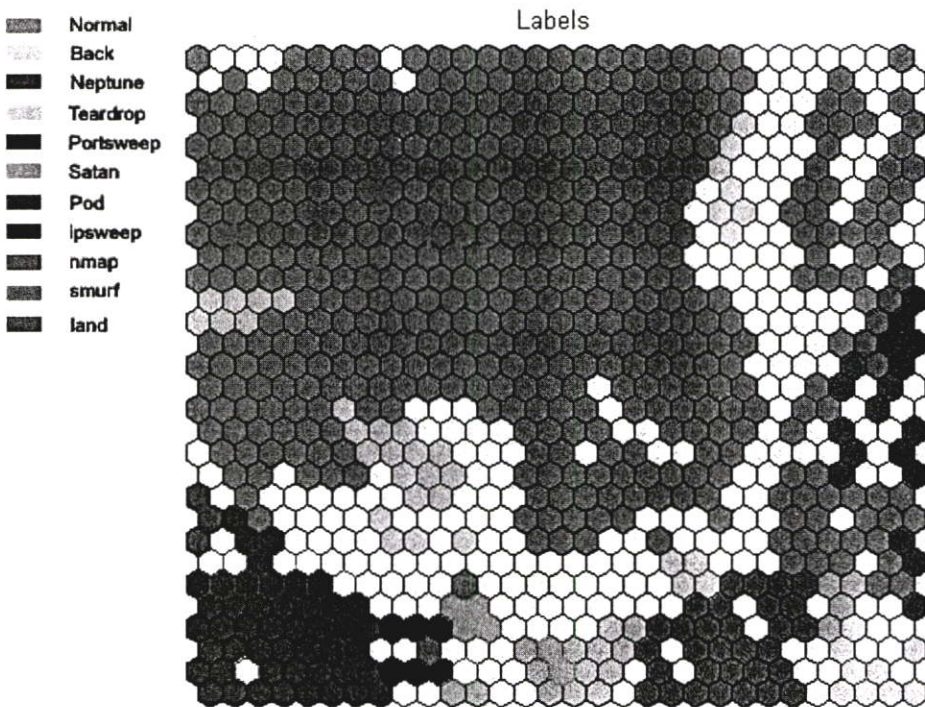
เมื่อทำการเตรียมชุดข้อมูลพฤติกรรมกรบุกรูกระบบเครือข่าย พร้อมกับค่าพารามิเตอร์ในการสอนระบบเริ่มต้นแล้ว ผู้วิจัยได้ทำการทดลองได้ผลดังตารางที่ 5.3 เพื่อทำการเปรียบเทียบเวลาในการสอนระบบ พร้อมกับตรวจสอบการซ้อนทับกันของกลุ่มข้อมูล

ตารางที่ 5.3 ผลการทดสอบ

ขนาด SOM	เวลา	จำนวนจัดกลุ่ม	จำนวนโหนด ตาย	จำนวนโหนดทับ ซ้อนสูงสุด
10x10	21 ชม. 10 นาที	8	39	5 ชนิด
15x15	21 ชม. 40 นาที	9	45	7 ชนิด
20x20	23 ชม. 20 นาที	9	69	6 ชนิด
25x25	24 ชม.10 นาที	10	151	5 ชนิด
28x28	24 ชม. 20 นาที	10	237	5 ชนิด
29x29	24 ชม. 30 นาที	10	238	5 ชนิด
30x30	24 ชม.30 นาที	11	254	5 ชนิด
35x35	24 ชม.50 นาที	11	368	5 ชนิด

40x40	25 ชม.5 นาที	11	480	5 ชนิด
45x45	25 ชม.15 นาที	11	608	5 ชนิด
50x50	25 ชม.45 นาที	11	750	5 ชนิด
55x55	26 ชม.15 นาที	11	908	5 ชนิด
60x60	26 ชม.35 นาที	11	1080	5 ชนิด

จากตารางที่ 5.3 การทดสอบได้ผลปรากฏว่า แผนภาพที่มีขนาด 30 x 30 หรือจำนวน 900 โหนด เป็นแผนภาพที่มีการกระจายกลุ่มพฤติกรรมการบุกรุกได้ดีกว่า แผนภาพขนาดอื่นๆ ดังรูปที่ 5.1 เป็นภาพแสดงการจัดกลุ่มของข้อมูลพฤติกรรมการบุกรุกที่สามารถแสดงกลุ่มข้อมูลได้ครบตามกลุ่มข้อมูลของพฤติกรรมชนิดต่างๆ ซึ่งมีทั้งหมด 11 กลุ่ม



รูปที่ 5.1 แสดงแผนภาพ SOM ขนาด 900 โหนด หรือขนาด 30 x 30

และเมื่อวัดประสิทธิภาพในการตรวจจับผู้บุกรุก detection rate ผลปรากฏว่าเปอร์เซ็นต์ความแม่นยำในการตรวจจับโดยรวมมีความแม่นยำสูงขึ้นมากกว่าเปอร์เซ็นต์ของแผนภาพในการทดลองที่ 1 ดังตารางที่ 5.4 ส่วนค่าความผิดพลาด false positive rate มีเปอร์เซ็นต์ที่ต่ำลงซึ่งแสดงให้เห็นว่าระบบการตรวจจับมีความผิดพลาดน้อยลง ดังตารางที่ 5.5

ตารางที่ 5.4 แสดงผลของการวัดค่า detection rate

Attack Type	detection rate
Normal	99.77%
Neptune	99.97%
Smurf	99.74%
Back	99.50%
Satan	93.79%
Ipsweep	93.10%
Portsweep	66.67%
Teardrop	98.96%
Pod	100.00%
Land	100.00%
Nmap	50.00%

ตารางที่ 5.5 แสดงผลของการวัดค่า false positive

SOM	false positive
ลำดับชั้นเดียว	0.23%

ในการทดลองที่ 1 นี้ผู้วิจัยนำเสนอการทดลองในส่วนเบื้องต้นเพื่อประยุกต์ SOM เข้ามาใช้ในการตรวจจับพฤติกรรมการบุกรุกระบบเครือข่าย โดยกำหนดขนาดของ SOM เป็นขนาดต่างๆ ดังตารางที่ 5.3 แล้วทำการทดสอบระบบเพื่อจับเวลาในการสอนระบบ ลำดับต่อไป คือ การจัดกลุ่มของแผนภาพ SOM ตามด้วยการนับจำนวนโหนดตาย และลำดับสุดท้าย คือ การตรวจสอบโหนดที่มีข้อมูลทับซ้อนว่าโหนดใดมีข้อมูลทับซ้อนมากที่สุดและมีจำนวนชนิดข้อมูลทับซ้อนสูงสุดกี่ชนิด เพื่อใช้ในการพิจารณาเลือกแผนภาพที่ดีที่สุด และมีประสิทธิภาพที่ดีที่สุดใน การตรวจจับการบุกรุกระบบเครือข่าย ซึ่งผลปรากฏว่าแผนภาพ SOM ที่เหมาะสม คือ แผนภาพ SOM ขนาด 30x30 หรือ จำนวน 900 โหนด แล้วนำไปทดสอบประสิทธิภาพของการตรวจจับการบุกรุกได้ผลดังตารางที่ 5.4 และ 5.5

ส่วนในการทดลองที่ 2 จะเป็นการทดลองที่ผู้วิจัยทำการทดลองเพื่อนำเสนอแนวคิดในการสร้างแผนภาพ SOM ขึ้นซ้อนกับแผนภาพ SOM ในชั้นแรก โดยคาดหวังว่าจะช่วยลดปัญหาการทับซ้อนของการจัดกลุ่มข้อมูลของแผนภาพ ดังรายละเอียดในหัวข้อถัดไป

5.2 การทดลองที่ 2 การตรวจจับการบุกรุกโดยใช้ SOM แบบหลายลำดับชั้น

5.2.1 จุดประสงค์ของการทดลอง

ทดสอบแนวคิดของผู้วิจัยในการสร้างแผนภาพ SOM ขึ้นซ้อนกับแผนภาพ SOM ในชั้นแรก คาดหวังช่วยลดปัญหาการทับซ้อนของการจัดกลุ่มข้อมูลพฤติกรรมกรการบุกรุกของแผนภาพ SOM โดยมีการกำหนดให้ค่าคงที่มีค่าเท่ากับ 0.8, 0.9 และ 1.0 เพื่อใช้ในการพิจารณาสร้างแผนภาพ SOM ชั้นใหม่ และนำไปวัดค่าประสิทธิภาพของการตรวจจับพฤติกรรมกรการบุกรุก

5.2.2 ขั้นตอนการทดลอง

ผู้วิจัยได้นำแผนภาพ SOM ลำดับชั้นเดียวในการทดลองที่ 1 มาประยุกต์สร้างแผนภาพ SOM เพิ่มเติม โดยกำหนดค่าคงที่ (β) เท่ากับ 0.8, 0.9 และ 1.0 เพื่อใช้ในการพิจารณาเพิ่มขึ้นของแผนภาพ SOM ซ้อนทับแผนภาพ SOM ในชั้นแรกและต่อไป โดยทำตามอัลกอริธึมในบทที่ 4 ในหัวข้อ 4.5 แล้วทำการสร้างแผนภาพ SOM ใหม่ขึ้นมาซ้อนทับแผนภาพ SOM เดิม ผลปรากฏดังตารางที่ 5.6 และเมื่อวัดประสิทธิภาพในการตรวจจับผู้บุกรุก detection rate และวัดค่าความผิดพลาด false positive rate ผลปรากฏดังตารางที่ 5.7 และตารางที่ 5.8

ตารางที่ 5.6 แสดงผลของจำนวนโหนดในแผนภาพ SOM และจำนวนชั้นของแผนภาพ SOM

ค่า β	จำนวนโหนดเป็น	จำนวนโหนดตาย	จำนวนโหนดทั้งหมด	จำนวนชั้นของแผนภาพ SOM
0.8	1,022	1,953	2,975	3
0.9	1,045	1,930	2,975	3
1.0	1,152	1,848	3,000	4

ตารางที่ 5.6 ผู้วิจัยได้นำเสนอผลของการทดลองในส่วนของจำนวนโหนดที่เป็นตัวแทนกลุ่มของข้อมูลพฤติกรรมกรการบุกรุก และจำนวนชั้นทั้งหมดของแผนภาพ SOM เพื่อใช้ในการพิจารณาว่า ค่าคงที่ (β) ที่กำหนดนั้นมีผลทำให้มีจำนวนโหนดเป็น ในแผนภาพ SOM จำนวนเท่าไร และสามารถสร้างแผนภาพได้จำนวนชั้นเท่าไร แล้วจึงนำไปทดสอบค่าความมีประสิทธิภาพต่อไป

ตารางที่ 5.7 แสดงผลของการวัดค่า detection rate

Attack Type	$\beta = 0.8$	$\beta = 0.9$	$\beta = 1.0$
Normal	99.84%	99.88%	99.90%
Neptune	100.00%	100.00%	100.00%
Smurf	99.86%	99.90%	99.90%

ตารางที่ 5.7 (ต่อ)

Back	100.00%	100.00%	100.00%
Satan	96.92%	96.95%	96.95%
Ipsweep	94.27%	97.47%	98.70%
PortswEEP	92.67%	94.53%	96.43%
Teardrop	100.00%	100.00%	100.00%
Pod	100.00%	100.00%	100.00%
Land	100.00%	100.00%	100.00%
Nmap	97.46%	98.32%	99.72%

ตารางที่ 5.8 แสดงผลของการวัดค่าfalse positive

SOM	$\beta = 0.8$	$\beta = 0.9$	$\beta = 1.0$
false positive	0.16%	0.12%	0.10%

จากการทดลองที่ 2 ผลปรากฏว่าการสร้างแผนภาพแบบมัลติเลเยอร์ SOM นั้นทำให้มีประสิทธิภาพในการตรวจจับผู้บุกรุกได้ดีกว่า แผนภาพ SOM แบบลำดับชั้นเดียว โดยดูได้จากตารางที่ 5.7 โดยเฉพาะในการตรวจจับพฤติกรรมการบุกรุกประเภท Nmap เพราะได้ค่าเปอร์เซ็นต์การตรวจจับที่สูงขึ้นถึง 97.46 เปอร์เซ็นต์ จากเดิม 50 เปอร์เซ็นต์ในการทดลองที่ 1 และค่าความผิดพลาดก็เพียง 0.16 เปอร์เซ็นต์ซึ่งน้อยกว่าจากเดิม 0.23 เปอร์เซ็นต์ในการทดลองที่ 1

ส่วนจำนวนโหนดเป็น และจำนวนชั้นทั้งหมดของแผนภาพที่ได้ ทำให้ผู้วิจัยสังเกตพบว่ามีความสัมพันธ์กับเวลาในการสอนระบบและทดสอบระบบเพราะว่า ยังมีจำนวนโหนดในการคำนวณน้อยก็ใช้เวลาน้อย และยังมีจำนวนชั้นน้อยเท่าไรก็ยังสามารถตรวจสอบพฤติกรรมได้เร็วมากขึ้นด้วยเช่นกัน แต่ก็ยังไม่สามารถสรุปได้ว่าถ้ากำหนดค่าคงที่เท่ากับ 0.8 แล้วจะทำให้ระบบตรวจจับการบุกรุกระบบเครื่องนี้มีประสิทธิภาพดีที่สุด ดังนั้นผู้วิจัยจึงได้ดำเนินการกำหนดให้ค่าคงที่มีค่าเป็น 0.9 และ 1.0 ตามลำดับ และผลปรากฏว่าเปอร์เซ็นต์ในการตรวจจับเมื่อกำหนดค่าคงที่ 1.0 นั้นมีเปอร์เซ็นต์การตรวจจับที่ดีที่สุด ส่วนเปอร์เซ็นต์ความผิดพลาดก็มีเปอร์เซ็นต์ความผิดพลาดก็ลดลงจากเดิม 0.16 เป็น 0.1 เปอร์เซ็นต์

5.3 การทดลองที่ 3 การทดสอบกับข้อมูลด้วย Back propagation เปรียบเทียบกับ SOM แบบหลายลำดับชั้น

5.3.1 จุดประสงค์ของการทดลอง

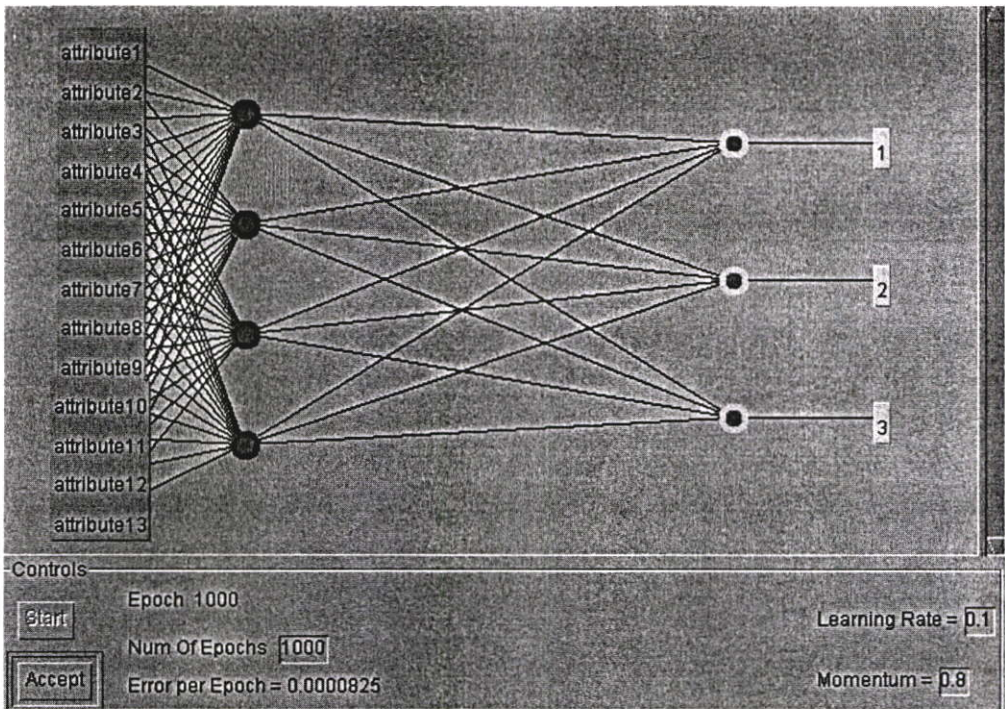
ทดสอบการจำแนกประเภทของข้อมูล ซึ่งประกอบไปด้วย Iris Dataset, Wine Dataset และพฤติกรรมกรนุกรุก KDD Cup จำนวน 1,000 เรคคอร์ด 5,000 เรคคอร์ด ตามลำดับด้วย Back propagation เพื่อนำผลที่ได้พิจารณาเปรียบเทียบกับโมเดล SOM แบบหลายลำดับชั้น

5.3.2 ขั้นตอนการทดลอง

ผู้วิจัยได้นำชุดข้อมูลมาตรฐาน Iris, Wine, KDD Cup 1,000, KDD Cup 5,000 ซึ่งแบ่งข้อมูลออกเป็น 20% สำหรับการฝึกสอนนิวรอลเน็ตเวิร์ค อีก 80 % สำหรับการทดสอบนิวรอลเน็ตเวิร์ค โดยมีการกำหนดค่าเริ่มต้นของ Back propagation คือ อินพุตเลเยอร์จำนวน 13 นิวรอล 1 เลเยอร์ ฮิดเดนเลเยอร์จำนวน 4 นิวรอล 1 เลเยอร์ เอาท์พุตเลเยอร์ 3 นิวรอล 1 เลเยอร์ สำหรับการทดลองกับชุดข้อมูล Wine และอินพุตเลเยอร์จำนวน 4 นิวรอล 1 เลเยอร์ ฮิดเดนเลเยอร์จำนวน 2 นิวรอล 1 เลเยอร์ เอาท์พุตเลเยอร์ 3 นิวรอล 1 เลเยอร์ สำหรับการทดลองกับชุดข้อมูล Iris และอินพุตเลเยอร์จำนวน 41 นิวรอล 1 เลเยอร์ ฮิดเดนเลเยอร์จำนวน 26 นิวรอล 1 เลเยอร์ เอาท์พุตเลเยอร์จำนวน 11 นิวรอล 1 เลเยอร์ สำหรับการทดลองกับชุดข้อมูล KDD Cup 1,000, KDD Cup 5,000 เรคคอร์ด ตามลำดับผลปรากฏดังตารางที่ 5.10

ตารางที่ 5.9 แสดงรายละเอียดชุดข้อมูลมาตรฐาน

Name	Instances	Attributes	Classes
Wine	178	13	3
Iris	150	4	3
KDD Cup	1,000	41	11
KDD Cup	5,000	41	11



รูปที่ 5.2 แสดงตัวอย่างผลการทดลองกับชุดข้อมูล Wine

ตารางที่ 5.10 แสดงผลของจำนวนการจำแนกประเภทข้อมูล

	Name	BP	MLSOM
Correct	Wine	92.2535 %	95.47 %
	Iris	95.467 %	96.375 %
	KDD Cup 1000	94.572 %	97.852 %
	KDD Cup 5000	92.367 %	98.526 %
Incorrect	Wine	7.7465 %	4.53 %
	Iris	4.533 %	3.625 %
	KDD Cup 1000	5.428 %	2.148 %
	KDD Cup 5000	7.633 %	1.474 %

จากตารางที่ 5.10 ผู้วิจัยได้นำเสนอการเปรียบเทียบการจำแนกชนิดของข้อมูล โดยวิธี Back propagation กับ SOM แบบหลายลำดับชั้น ผลปรากฏว่าการจำแนกประเภทแบบ SOM แบบหลายลำดับชั้นมีผลการจำแนกประเภทที่ถูกต้องมากกว่า และมีผลการจำแนกผิดพลาดน้อยกว่า

จากการทดลองที่ 3 ผลปรากฏว่าการจำแนกข้อมูลแบบ Back propagation มีผลการจำแนกประเภทข้อมูลที่ถูกต้องน้อยกว่า และมีผลการจำแนกผิดพลาดมากกว่า โมเดล SOM แบบหลายลำดับชั้น ดังนั้นจากการทดลองนี้ โมเดล SOM แบบหลายลำดับชั้นจึงมีแนวโน้มในการ

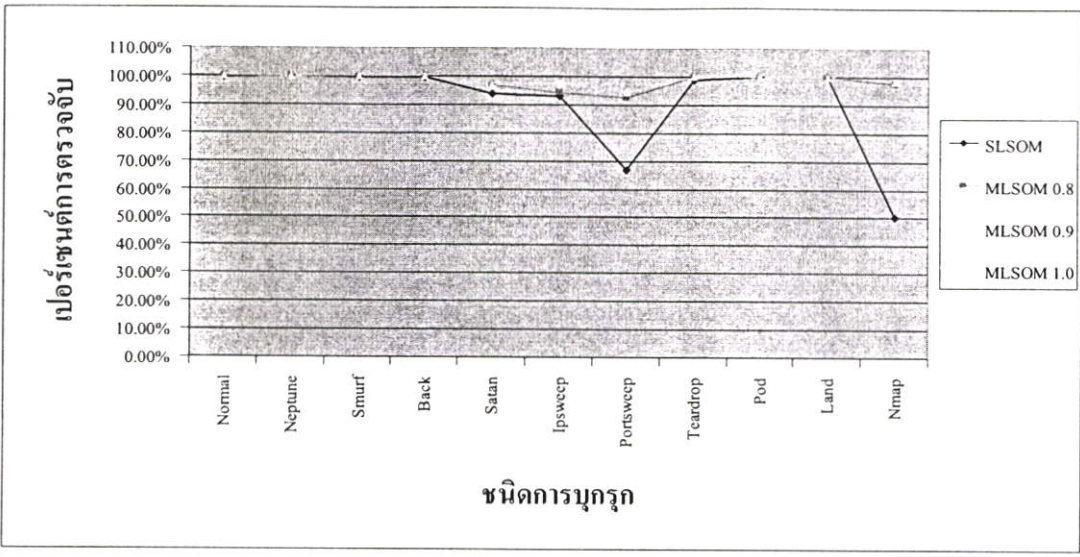
จำแนกประเภทได้ดีกว่า Back propagation และจากการทดลองที่ 2 ผลปรากฏว่าการสร้างแผนภาพแบบมัลติเลเยอร์ SOM นั้นทำให้มีประสิทธิภาพในการตรวจจับผู้บุกรุกได้ดีกว่าแผนภาพ SOM แบบลำดับชั้นเดียว โดยดูได้จากตารางที่ 5.7 โดยเฉพาะในการตรวจจับพฤติกรรมการบุกรุกประเภท Nmap เพราะได้ค่าเปอร์เซ็นต์การตรวจจับที่สูงขึ้นถึง 97.46 เปอร์เซ็นต์ จากเดิม 50 เปอร์เซ็นต์ในการทดลองที่ 1 และค่าความผิดพลาดก็เพียง 0.16 เปอร์เซ็นต์ซึ่งน้อยกว่าจากเดิม 0.23 เปอร์เซ็นต์ในการทดลองที่ 1

ส่วนจำนวน โหนดเป็น และจำนวนชั้นทั้งหมดของแผนภาพที่ได้ ทำให้ผู้วิจัยสังเกตพบว่ามีความสัมพันธ์กับเวลาในการสอนระบบและทดสอบระบบเพราะว่า ยังมีจำนวนโหนดในการคำนวณน้อยก็ใช้เวลาน้อย และยังมีจำนวนชั้นน้อยเท่าไรก็ยังสามารถตรวจสอบพฤติกรรมได้เร็วมากขึ้นด้วยเช่นกัน แต่ก็ยังไม่สามารถสรุปได้ว่าถ้ากำหนดค่าคงที่เท่ากับ 0.8 แล้วจะทำให้ระบบตรวจจับการบุกรุกระบบเครื่องนี้มีประสิทธิภาพดีที่สุด ดังนั้นผู้วิจัยจึงได้ดำเนินการกำหนดให้ค่าคงที่มีค่าเป็น 0.9 และ 1.0 ตามลำดับ และผลปรากฏว่าเปอร์เซ็นต์ในการตรวจจับเมื่อกำหนดค่าคงที่ 1.0 นั้นมีเปอร์เซ็นต์การตรวจจับที่ดีที่สุด ส่วนเปอร์เซ็นต์ความผิดพลาดก็มีเปอร์เซ็นต์ความผิดพลาดก็ลดลงจากเดิม 0.16 เป็น 0.1 เปอร์เซ็นต์

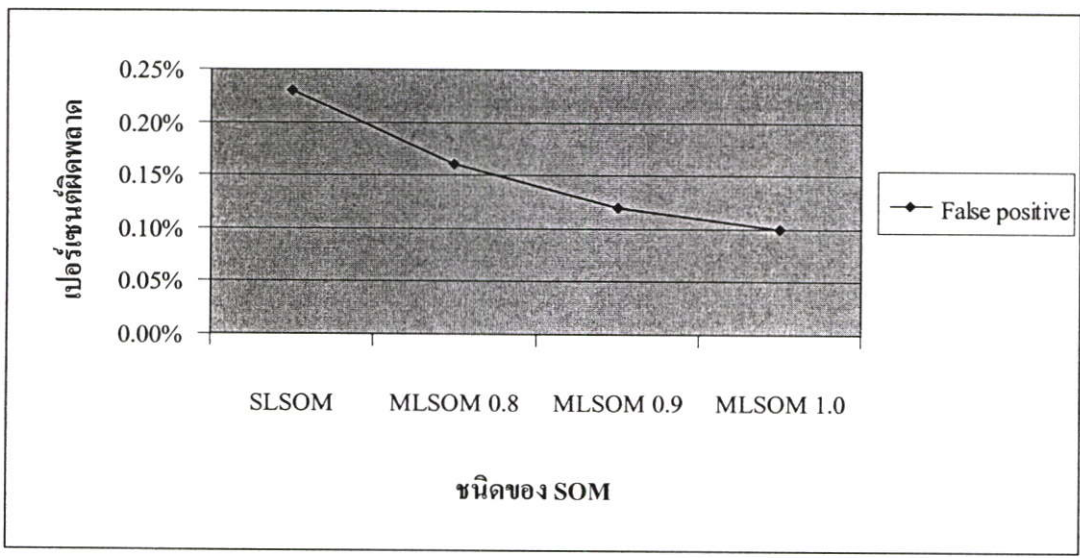
5.3 สรุป

จากการทดลองผลปรากฏว่า ในการสร้างแผนภาพ SOM แบบหลายลำดับชั้นนั้นมีผลเปอร์เซ็นต์การตรวจจับที่ดีกว่า ดังรูปที่ 5.2 จะเห็นได้ว่าเมื่อมีการกำหนดค่า β เท่ากับ 1.0 แล้วมีผลทำให้เปอร์เซ็นต์ในการตรวจจับผู้บุกรุกดีกว่าเมื่อมีการกำหนดค่า β เท่ากับ 0.8, 0.9 และแผนภาพ SOM แบบลำดับชั้นเดียวในการตรวจจับผู้บุกรุก และเมื่อพิจารณารูปที่ 5.3 ก็สามารถแสดงให้เห็นถึงประสิทธิภาพในการตรวจจับผู้บุกรุก เพราะว่ามีเปอร์เซ็นต์ในการตรวจจับผู้บุกรุกผิดพลาดต่ำลงอย่างเห็นได้ชัด ทำให้ระบบการตรวจจับผู้บุกรุกนี้จึงมีความน่าเชื่อถือมากยิ่งขึ้น

และในบทถัดไปจะนำเสนอถึงบทสรุปผลงานวิจัย และนำเสนอข้อเสนอแนะต่างๆ เพิ่มเติม



รูปที่ 5.2 แสดงเปอร์เซ็นต์การตรวจจับ Detection rate



รูปที่ 5.3 แสดงเปอร์เซ็นต์ False positive

บทที่ 6

สรุปผลการวิจัย และข้อเสนอแนะ

6.1 สรุปผลการวิจัย

ผู้วิจัยได้นำเสนอแนวคิดในการประยุกต์ใช้ SOM ในการตรวจจับพฤติกรรมการบุกรุกระบบเครือข่าย โดยในเบื้องต้นเกิดแนวความคิดในการพยายามนำเสนอรูปแบบพฤติกรรมการบุกรุกระบบเครือข่ายคอมพิวเตอร์ให้สามารถจำแนกและแยกประเภทพฤติกรรมการบุกรุกออกมาได้อย่างชัดเจน เพื่อเป็นประโยชน์ต่อผู้ดูแลระบบเครือข่ายคอมพิวเตอร์ โดย SOM มีลักษณะเด่น คือสามารถนำเสนอแผนภาพการจัดกลุ่มของข้อมูล ทำให้สามารถเห็นภาพกลุ่มของข้อมูลได้ และในงานวิจัยนี้ผู้วิจัยได้พัฒนา SOM แบบหลายลำดับชั้น โดยมีจุดประสงค์เพื่อเพิ่มประสิทธิภาพในการตรวจจับผู้บุกรุกให้มีประสิทธิภาพมากยิ่งขึ้น เนื่องจากในการตรวจจับผู้บุกรุกโดยใช้ SOM แบบลำดับชั้นเดียวนั้นมีปัญหาเรื่องการซ้อนทับกับของพฤติกรรมการบุกรุกที่แสดงในแผนภาพ SOM อันเนื่องมาจากโหนดตัวแทนกลุ่มข้อมูลพฤติกรรมของแผนภาพ แสดงตัวแทนชนิดของพฤติกรรมการบุกรุกมากกว่าสองชนิดในโหนดเดียวกัน จึงเป็นผลให้ระบบไม่สามารถจำแนกประเภทการบุกรุกได้อย่างมีประสิทธิภาพ และในระบบการตรวจจับผู้บุกรุกโดยใช้ SOM แบบหลายลำดับชั้นนั้นได้แก้ไขปัญหาดังกล่าว โดยมีการสร้างอัตราส่วนของชนิดข้อมูลที่ตกในโหนดเดียวกันแล้วเปรียบเทียบค่าที่สูงที่สุด เพื่อนำมาพิจารณาว่าควรสร้างชั้นของแผนภาพ SOM เพิ่มหรือไม่ โดยใช้ข้อมูลที่อยู่ในโหนดที่มีข้อมูลที่ซ้อนทับกัน

ในการทดลองที่ 1 นี้ผู้วิจัยจึงทำการทดสอบหาขนาดที่เหมาะสมของแผนภาพ SOM ที่มีประสิทธิภาพในการตรวจจับพฤติกรรมการบุกรุก ผลปรากฏว่าแผนภาพที่เหมาะสมนั้นมีขนาด 30x30 หรือ จำนวนโหนด 900 โหนด เนื่องจากว่าสามารถจัดกลุ่มได้เท่ากับกลุ่มของข้อมูลพฤติกรรมการบุกรุกตัวอย่างที่นำมาทดสอบ และใช้เวลาน้อยที่สุด แต่ปัญหาของแผนภาพ SOM ลำดับชั้นเดียว คือ กลุ่มพฤติกรรมการบุกรุกบางชนิดมีการทับซ้อนกันในแผนภาพของ SOM ลำดับชั้นเดียว ส่งผลให้ประสิทธิภาพของระบบการตรวจจับพฤติกรรมการบุกรุกขาดประสิทธิภาพไป

และในปัญหานี้เองผู้วิจัยจึงนำเสนอการแก้ไขโดยการสร้างมัลติเลเยอร์ SOM โดยคาดหวังว่าจะสามารถช่วยแก้ปัญหาดังกล่าวได้ จึงได้นำเสนอการทดลองที่ 2 ผลปรากฏว่าสามารถแก้ปัญหาในเรื่องของการซ้อนทับกันของข้อมูลได้ และประสิทธิภาพของระบบการตรวจจับโดยรวมมีประสิทธิภาพในการตรวจจับพฤติกรรมการบุกรุกได้แม่นยำ พร้อมกับมีความผิดพลาดต่ำลง เนื่องจากข้อมูลพฤติกรรมการบุกรุกที่ซ้อนทับกันนั้น ได้ถูกกระจายออกไปยังแผนภาพ SOM ใหม่ ซึ่งทำให้ช่วยลดปัญหาเรื่องการซ้อนทับของข้อมูลได้ และเมื่อไม่มีการซ้อนทับของข้อมูล

แล้วประสิทธิภาพในการตรวจจับก็ดีขึ้นตามไปด้วย ผลสุดท้ายจึงทำให้ระบบมีประสิทธิภาพโดยรวมดีขึ้น และมีความน่าเชื่อถือมากขึ้น

6.1 ข้อเสนอแนะ

ประการแรกเนื่องจากข้อมูลที่ใช้ในการทดสอบระบบเป็นข้อมูลที่จำลองขึ้นของ KDD Cup 1999 ซึ่งเป็นข้อมูลมาตรฐานที่เป็นที่นิยมสำหรับนักวิจัยทั่วไป ซึ่งในโอกาสต่อไปผู้วิจัยจะค้นหาข้อมูลมาตรฐานใหม่ๆ ที่เป็นที่ยอมรับสำหรับนักวิจัยเพื่อทำการทดสอบต่อไป

ประการที่สองข้อมูลที่ทำการทดลองมีบางคุณลักษณะในข้อมูลชุดนี้ซึ่งไม่มีค่าเปลี่ยนแปลงใดๆ เมื่อพิจารณาแล้วอาจสามารถตัดออกไปจากฐานข้อมูลได้ ก็จะช่วยให้การประมวลผลของ SOM ทำได้เร็วมากขึ้น ซึ่งในโอกาสต่อไปผู้วิจัยจะทดสอบระบบกับข้อมูลที่มีการพิจารณาเลือกคุณลักษณะออกเป็นกลุ่มต่างๆ เพื่อทดสอบเปรียบเทียบต่อไป

ประการที่สามคือโมเดลที่นำเสนอในงานวิจัยนี้เป็นการนำเสนอแนวคิดในการแก้ปัญหาการซ้อนทับของพฤติกรรมการบุกรุกของแผนภาพ SOM แบบลำดับชั้นเดียว ซึ่งผลการทดลองขั้นต้นแสดงให้เห็นว่าโมเดลที่นำเสนอ ให้ประสิทธิภาพในการจำแนกประเภทการบุกรุกที่ดีกว่า ในโอกาสต่อไปผู้วิจัยจะวิจัยถึงผลกระทบของรูปแบบการรวมกลุ่ม โหนดและจำนวนสมาชิกต่อประสิทธิภาพของโมเดล

บรรณานุกรม

- [1] S. A. Hofmeyr, "An Immunological Model of Distributed Detection and Its Application to Computer Security," PhD thesis, University of New Mexico, Albuquerque, New Mexico, 1999.
- [2] W. Lee, R. A. Nimbalkar, K. Yee and S. J. Stolfo, "A Data Mining and CIDF Based Approach for Detecting Novel and Distributed Intrusions," In Proceedings of the 3rd International Workshop on Recent Advances in Intrusion Detection (RAID 2000), October 2000.
- [3] W. Lee and S. J. Stolfo, "Data mining approaches for intrusion detection," In Proceedings of the 7th USENIX Security Symposium, San Antonio, TX, 1998.
- [4] P. Lichodziejewski, A. n. Zincir-Heywood and M. I. Heywood, "Host Based Intrusion Detection Using Self-Organizing Maps," In Proceedings of the 2002 IEEE World Congress on Computational Intelligence, 2002.
- [5] C. Jirapummin, N. Wattanapongsakorn, and P. Kanthamanon, "Hybrid Neural Networks for Intrusion Detection Systems," [Online], Available: http://dbvis.fmi.uni-konstanz.de/members/panse/seminar_ws0203/
- [6] H. Kayacik, A. Zincir-Heywood, and M. Heywood, "On the capability of an SOM based intrusion detection system," In Proceedings IEEE Int. Joint Conf. Neural Networks (IJCNN'03), pp. 1808-1813, 2003.
- [7] D. Denning, "An intrusion-detection model," IEEE Trans. Software Eng., vol. SE-13, no. 2, pp. 222-232, Feb 1987.
- [8] Helge Ritter, Thomas Martinetz and Klaus Schulten, "Neural computation and self-organizing maps : an introduction," Massachusetts : Addison-Wesley, 1992.
- [9] Xia Lin, Dagobert Soergal, Gary Marchioninl, "A Self-Organizing Semantic Map," ACM, 1991.
- [10] Kaski S., Honkela T., Lagus K., and Kohonen T., "WEBSOM-self-organizing maps of document collections," Neurocomputing, volume 21, 1998, pp. 101-117.
- [11] Kohonen T., "Self-organization of very large document collections: State of the art," ICANN98, Springer, London, 1998, pp. 65-74.
- [12] A.Rauber, D. Merkl., "The SOMLib Digital Librery System," Proceedings of the 3rd Europ. Conf. on Research and Advanced Techonology for Digital Libraries (ECDL'99), Paris, France, September 22-24 1999.

- [13] R.Baeza-Yates and B. Ribeiro-Neto., "Modern Information Retrieval," New York: ACM-Press. 1999.
- [14] Qing Ma, Min Zhang, Ming Zhou., "Self-Organization of Chinese Semantic Maps Using TFIDF Term Weighting," The Second Workshop on Natural Language Processing and Neural Networks, Tokyo, Japan, November, 2001.
- [15] K.Fox, R. Henning, and J. Reed, "A neural network approach toward intrusion detection," In Proc. 13th Nat. Computer Security Conf., Washington, DC, 1990.
- [16] J. Cannady, "Artificial neural networks for misuse detection," In Proceedings of the 1998 National Information Systems Security Conference (NISSC'98) October 5-8 1998. Arlington, VA., page 443-456, 1998.
- [17] S.Stolfo et al., The Third International Knowledge Discovery and Data Mining Tools Competition., <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>

ภาคผนวก



ISSN 0125-1724

วิศวกรรม

ลาดกระบัง

คณะวิศวกรรมศาสตร์ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

LADKRABANG ENGINEERING JOURNAL

ปีที่ 22 ฉบับที่ 4

ธันวาคม 2548

1.	Self-Organizing Map กับการจัดอันดับส่วนที่การตรวจงานการปลูก สุรพล วิชาประคิษฐ์ เอื้อน บินเงิน	1
2.	การวัดคุณสมบัติใน Self-Organizing Map โดยใช้เจเน็ตอัลกอริทึม กสวนดี ศรีกุดนาค พุทธิเทพ โรงนวนตุ ไพฑูรย์ ศรีนิล เอื้อน บินเงิน	7
3.	การวิเคราะห์โครงสร้างการเชื่อมต่อบน IPv6 ภายในองค์กรและการประยุกต์ใช้งาน สุวิยา เจริญสุดิดาวร กอบชัย เศรษฐนาถ	13
4.	การออกแบบและสร้างหม้อแปลงไฟฟ้าแรงสูงความถี่สูงขนาด 20 kV 2 MVA กิตติพงษ์ ตันมิตร อัญญา สุขศรี ชัยพร ชัดโคตร	19
5.	การออกแบบและวิเคราะห์วงจรเรียงกระแสเซมิคอนดักเตอร์ที่มีวงจรปรับรูปคลื่นแรงดันเอาต์พุต และกระแสอินพุต สกต กสิณรัตน์ วิจิตร กิณเรศ	25
6.	ผลกระทบของน้ำยาเคมี และน้ำ Di ที่มีต่อคุณสมบัติทางไฟฟ้าและแม่เหล็กของหัวขานเขียนข้อมูล สมเกียรติ ปราบก วิสุทธิ วิสุทธิเรือง สัตวาลัย สุภาดิ	31
7.	การสังเคราะห์คาร์บอนนาโนทิวป์ด้วยวิธี CVD แบบลดความดันร้อนที่ความดัน 1 บรรยากาศ โดยใช้แอลกอฮอล์ และไนโตรเจนเป็นก๊าซพาหะ ณธวรรษ กลิกรุ่งโรจน์ ปฎิคม ศรีหมพล สุวิชัย ชัยสิทธิ์ศักดิ์	36
8.	วงจรกำเนิดฟังก์ชันเชิงรีโหนดงานในโหมดกระแสใช้แรงดันต่ำด้วยเทคโนโลยี CMOS มนตรี คำเงิน วิฑูรช บุญนา กอบชัย เศรษฐนาถ	42
9.	วงจรบวกทางเวกเตอร์หลายหน้าที่ด้วย CMOS มนตรี คำเงิน คมกฤษ โภมเจตศิริ กอบชัย เศรษฐนาถ	46
10.	ผลกระทบของเครือข่ายที่ใช้ควบคุมหลักสั่งการควบคุมขยับเพียงชุดควบคุมเดียวของระบบแรงเหวี่ยงลงใหม่ วิจิตร แก้วโพธิ์เทียม กอบชัย เศรษฐนาถ	52
11.	การวิเคราะห์สมรรถนะของระบบ DS-QPSK CDMA โดยใช้ข้อสังเกตจากการางแบบนาคากามิ เกียรติวุฒิ จรภักดี กอบชัย เศรษฐนาถ	57

Self-Organizing Map หลายลำดับชั้นสำหรับการตรวจจับ การบุกรุก

Multi-layer Self-Organizing Map for Intrusion Detection

ศุรพล โรจนประดิษฐ์ เอียน ปิ่นเงิน
ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์
สำนักวิจัยการสื่อสารและเทคโนโลยีสารสนเทศ
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

บทคัดย่อ

การตรวจจับการบุกรุกเครือข่ายโดยอาศัยการแยกประเภทการบุกรุกเครือข่ายที่ใช้ Self-Organizing Map (SOM) เป็นการจัดข้อมูลพฤติกรรมการบุกรุกที่ได้จากคุณลักษณะทางระบบเครือข่ายหลายมิติ โดยแปลงให้อยู่ในรูปของแผนภาพ SOM สองมิติ ซึ่งจะให้ข้อมูลการบุกรุกระบบเครือข่ายที่มีลักษณะคล้ายกันถูกจัดกลุ่มให้อยู่โหนดเดียวกัน แต่ในบางกรณีจะเกิดการซ้อนทับกันของข้อมูล กล่าวคือข้อมูลโหนดเดียวกันอาจเกิดจากการบุกรุกเครือข่ายที่ต่างประเภทกัน ทำให้ไม่สามารถระบุประเภทของการบุกรุกระบบเครือข่ายได้อย่างชัดเจน ปัญหาหลักของแผนภาพที่มีข้อมูลประเภทการบุกรุกที่ทับซ้อนกันนั้น SOM แบบลำดับชั้นเดียวไม่สามารถแก้ปัญหาได้ ดังนั้นงานวิจัยนี้ได้นำเสนอการแยกประเภทการบุกรุกด้วยวิธี SOM แบบหลายลำดับชั้น โดยแยกการทำงานออกเป็นสองขั้นตอน ขั้นตอนแรกเป็นการแยกประเภทการบุกรุกเครือข่ายจากข้อมูลขั้นต้นด้วยวิธี SOM แบบลำดับชั้นเดียว และขั้นที่สองตรวจสอบ SOM ที่ได้ว่าโหนดใด มีอัตราการซ้อนทับของข้อมูลมากเกินกว่าที่กำหนดก็จะทำการแยกข้อมูลเฉพาะโหนดนั้นออกเป็นอีกหนึ่งลำดับชั้น จากการทดลองแสดงให้เห็นว่าด้วยวิธีการที่นำเสนอให้ผลค่าเปอร์เซ็นต์ Detection rate และค่าเปอร์เซ็นต์ False positive ดีกว่าการใช้ SOM แบบลำดับชั้นเดียว

Abstract

The classification of network intrusion detection by using Self-Organizing Map (SOM) is an intrusion behavior management that mapped the multi-dimensional features into two dimensional SOM. Consequently, the similarities of intrusion data are classified in the same node. However, there are overlapping data in some case such as data in a specific cluster have different types of attack. Therefore, it is difficult to identify the type of intrusion behavior. The problem of overlapping data can not be solved by single-layer SOM. This research presented the classification of network intrusion detection using multi-layer SOM in order to classify the types of intrusion. The process consists of two steps. First, the algorithm uses single-layer SOM to classify types of intrusion from the primary data. Second, the results from SOM are examined to determine the nodes that have overlapping rate exceeding the threshold. Then the data in one layer is distributed again. From the experiments, we found that the percentage of detection rate and false positive rate were better than single-layer SOM.

1. บทนำ

การตรวจจัดการบุกรุกเป็นส่วนสำคัญส่วนหนึ่งในการรักษาความปลอดภัยบนระบบเครือข่ายคอมพิวเตอร์ แต่เนื่องจากในปัจจุบันการพัฒนาทางด้านระบบเครือข่ายคอมพิวเตอร์ได้มีการพัฒนาไปอย่างรวดเร็วและมีการนำเสนอต่อสาธารณะมากขึ้น ก็ย่อมเกิดจำนวนรูปแบบการบุกรุกมากขึ้นตามไปด้วย ดังนั้นความสนใจในการพัฒนาระบบการตรวจจับผู้บุกรุกจึงมุ่งไปที่การใช้เทคนิคคาน่าโมนิงในการตรวจจัดการบุกรุกเครือข่าย [1] คือ Misuse detection และ Anomaly detection [2] เทคนิคแบบ Misuse detection นั้นจะมีลักษณะของการเก็บรูปแบบการบุกรุกไว้ในรูปแบบของฐานข้อมูลกฎ ซึ่งกระบวนการนี้จะได้ผลดีมาก ถ้าหากว่ามีการบุกรุกที่เคยมูลกรระบบเครือข่ายนี้มาก่อนและถูกจัดเก็บรูปแบบไว้ในฐานข้อมูลกฎแล้ว แต่ข้อเสียของกระบวนการนี้คือ ไม่สามารถตรวจจัดการบุกรุกระบบเครือข่ายรูปแบบใหม่ หรือว่าไม่เคยถูกจัดเก็บไว้ในฐานข้อมูลกฎได้ ส่วน Anomaly detection นั้นจะเป็นลักษณะการแยกประเภทพฤติกรรมการใช้งานระบบเครือข่ายของผู้ใช้งานปกติ (normal) ออกจากประเภทพฤติกรรมการใช้งานระบบเครือข่ายของผู้รบกวนต่อระบบเครือข่าย (abnormal) ซึ่งระบบที่ใช้กระบวนการนี้ในการตรวจจับก็จะทำการแจ้งเตือนว่าเป็นการบุกรุก ซึ่งข้อดีคือสามารถตรวจจับพฤติกรรมกรบุกรุกระบบเครือข่ายคอมพิวเตอร์ประเภทใหม่ๆ ได้ เทคนิคหนึ่งที่ถูกนำมาใช้คือนิวรอลเน็ตเวิร์ก [3,4]

SOM เป็นนิวรอลเน็ตเวิร์กแบบไม่มีผู้สอน ถูกนำมาประยุกต์ใช้ในการตรวจจัดการบุกรุก เช่นงานวิจัยของ H. Kayacik [5] ได้นำเสนอลักษณะของ SOM ตาม ลำดับขั้นเพื่อทำการจำแนกประเภทการบุกรุกระบบเครือข่ายโดยใช้คุณลักษณะเพียง 6 คุณลักษณะพื้นฐาน [2] ในลำดับขั้นที่หนึ่งจากนั้นนำโหนดที่ชนะที่ได้จากลำดับขั้นที่หนึ่ง ส่งไปทำการสร้าง SOM ในลำดับขั้นที่สอง และสุดท้ายนำโหนดชนะที่ได้ส่งต่อไปสร้าง SOM ในลำดับขั้นที่สาม ได้ค่า False positive rate 4.6% และค่า Detection rate 89% ส่วนงานวิจัยของ Susecla T. [6] ได้นำเสนอ SOM แบบตามลำดับขั้นเช่นเดียวกับงานวิจัยของ H. Kayacik ในการ

ตรวจจัดการบุกรุกแบบ Anomaly detection เช่นกัน แต่เป็นการคัดเลือกคุณลักษณะออกเป็น ห้ากลุ่มในการสร้าง SOM แต่ละลำดับขั้น ทั้งหมดสามลำดับขั้นและทำการเปรียบเทียบประสิทธิภาพในการตรวจจัดการบุกรุกระบบเครือข่ายในแต่ละกลุ่ม ได้ค่า False positive rate 0.34% และค่า Detection rate 99.63%

งานวิจัยนี้นำเสนอวิธีการตรวจจัดการบุกรุกโดยอาศัยการแยกประเภทข้อมูลการบุกรุกระบบเครือข่ายคอมพิวเตอร์โดยแสดงให้อยู่ในรูปแบบของ SOM แบบหลายลำดับขั้น โดยมีการคำนวณค่าอัตราส่วนค่าหนึ่งที่ได้จากการพิจารณาโหนดที่มีข้อมูลที่ทับซ้อนกันเพื่อใช้เป็นข้อกำหนดในการเพิ่มลำดับขั้น และทำการทดลองบนพื้นฐานของการบุกรุกระบบเครือข่ายแบบ Denial of Service กับ Probing และพฤติกรรมปกติ

2. Self-Organizing Map (SOM)

SOM [7] คือนิวรอลเน็ตเวิร์กแบบไม่มีผู้สอนที่ได้รับความนิยมสูง โมเดลประกอบไปด้วยโหนดของนิวรอลโหนดที่ในแต่ละโหนด i จะเป็นค่าเวกเตอร์น้ำหนักแทนด้วย $\vec{m}_i = (w_1, w_2, \dots, w_n)$ เมื่อ $w_j \in \mathcal{R}$ และ n คือจำนวนคุณลักษณะของอินพุตเวกเตอร์

กระบวนการเรียนรู้ของ SOM มีขั้นตอนดังนี้

1. เลือกอินพุตเวกเตอร์แบบสุ่มจากอินพุตโดเมน
2. นำอินพุตเวกเตอร์ $\vec{x}_i(t)$ ไปเปรียบเทียบกับเวกเตอร์ $\vec{m}_i(t)$ ของทุกๆ โหนดเพื่อหาโหนดชนะจากโหนดทั้งหมด
3. ปรับเวกเตอร์น้ำหนักของโหนดชนะ เพื่อให้โหนดชนะเข้าใกล้อินพุตมากขึ้น
4. ปรับเวกเตอร์น้ำหนักของโหนดใกล้เคียง (neighborhood nodes) เพื่อให้อินพุตเวกเตอร์ถัดไปที่มีค่าใกล้เคียงมีโหนดชนะใหม่อยู่ใกล้กัน
5. ทำซ้ำกระบวนการ 1 ถึง 4 จนกว่าจะถึงจำนวนรอบหรือเงื่อนไขที่กำหนด

โดยเงื่อนไขที่กำหนดนั้นอาจเป็นค่าความผิดพลาดโดยรวมของระบบเป็นไปตามที่กำหนด หรือค่าเวกเตอร์น้ำหนักของโหนดชนะในรอบก่อนหน้ามีค่าคงที่เป็นต้น

ในการหาโหนดชนะฟังก์ชันที่นิยมใช้ในการเปรียบเทียบคือฟังก์ชันระยะห่างเชิงยูคลิดีอัน (Euclidean distance) โดยนำมาคำนวณหาโหนดชนะ c ดังสมการที่ 1

$$c : \bar{m}_c(t) = \min_i \|\bar{x}(t) - \bar{m}_i(t)\| \quad (1)$$

สมการปรับค่าของโหนดชนะและโหนดใกล้เคียง

สามารถแสดงได้ดังสมการที่ 2

$$\bar{m}_i(t+1) = \bar{m}_i(t) + \alpha(t) \times h_{ci}(t) \times [\bar{x}(t) - \bar{m}_i(t)] \quad (2)$$

เมื่อ

t คือรอบที่ ในการเรียนรู้

$\bar{x}(t)$ คืออินพุตเวกเตอร์ ณ รอบที่ t

$\bar{m}_i(t)$ คือเวกเตอร์น้ำหนัก

$\alpha(t)$ คืออัตราการเรียนรู้ในรอบที่ t ซึ่งแสดงได้ดังสมการที่ 3

$$\alpha(t) = \alpha(0) \times \frac{T-t}{T} \quad (3)$$

เมื่อ

T คือจำนวนรอบทั้งหมด

t คือรอบที่ t

$h_{ci}(t)$ คือฟังก์ชันที่ใช้ในการกำหนดน้ำหนักในการปรับค่าโหนดใกล้เคียง ซึ่งโดยทั่วไปแล้วจะใช้ฟังก์ชันเกาส์เซียน (Gaussian Function) ดังสมการที่ 4

$$h_{ci}(t) = \exp\left(-\frac{\|c_i - r_j\|^2}{2\sigma^2(t)}\right) \quad (4)$$

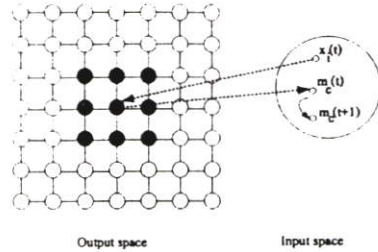
เมื่อ

$\|c_i - r_j\|^2$ แทน ระยะห่างของตำแหน่งของโหนด i กับโหนดชนะ c

$\sigma(t)$ แทน รัศมีของบริเวณโหนดใกล้เคียง

โดยปกติรัศมีของโหนดใกล้เคียงจะค่อยๆ ลดลงตามจำนวนรอบในการเรียนรู้ ดังสมการที่ 5

$$\sigma(t+1) = 1 + (\sigma(t) - 1) \times \frac{T-t}{T} \quad (5)$$



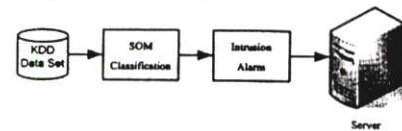
รูปที่ 1 โครงสร้างของ SOM และการเรียนรู้

ตัวอย่างจากรูปที่ 1 แสดง SOM ทรงกลมที่มีขนาด 7x7 อินพุตเวกเตอร์ $\bar{x}(t)$ ถูกเลือกแบบสุ่มจากโคมนของอินพุต ขั้นตอนไปหาโหนดชนะจากโหนดทั้งหมด จากรูปโหนดที่สี่เข้มสุดคือโหนดชนะ หลังจากที่ได้โหนดชนะ c ที่มีระยะห่างยูคลิดีอันน้อยที่สุดแล้วทำการปรับค่าเวกเตอร์น้ำหนัก $\bar{m}_c(t)$ เป็น $\bar{m}_c(t+1)$ ซึ่งจะเข้าใกล้อินพุตเวกเตอร์มากขึ้น หลังจากนั้นจะทำการปรับโหนดใกล้เคียงจากรูปคือโหนดต่างๆที่สีจางลงมา ความเข้มสีของโหนดแสดงให้เห็นถึงน้ำหนักของการปรับค่าโหนดโดยที่สี่เข้มหรือโหนดชนะจะมีการปรับค่าน้ำหนักมากกว่าโหนดสีจางหรือโหนดใกล้เคียง

3. ขั้นตอนการทำงานของระบบและการประยุกต์ใช้ SOM

3.1. ขั้นตอนการทำงานของระบบ

ขั้นตอนการทำงานของระบบที่เรานำ SOM เข้ามาประยุกต์ใช้นั้นดังแสดงในรูปที่ 2



รูปที่ 2 แสดงขั้นตอนการทำงานของระบบ

ในส่วนตัวแรกคือข้อมูลที่ใช้ในงานวิจัยนี้เป็นข้อมูลที่ได้จาก KDD cup 1999 [8] เป็นข้อมูลที่จำลองพฤติกรรมการบุกรุกใน 4 ประเภทหลักดังนี้

- Denial of Service ผู้บุกรุกพยายามโจมตีให้ระบบหยุดให้บริการ เช่น smurf

- Probing ผู้บุกรุกพยายามตรวจสอบหาจุดอ่อนของระบบ เช่น portsweep
- R2L ผู้บุกรุกไม่มียูสเซอร์ในระบบแต่พยายามเจาะระบบ เช่น guess password
- U2R ผู้บุกรุกพยายามเข้าสู่ระบบโดยใช้สิทธิ์ของซูเปอร์ยูสเซอร์ เช่น buffer overflow

และมีคุณลักษณะ (Feature) 41 คุณลักษณะได้มาจากการคักจับข้อมูลที่มีการสื่อสารกันในระบบเครือข่ายคอมพิวเตอร์แล้วนำมาใช้เทคนิคคาค่าไบนารีเปลี่ยนให้อยู่ในรูปของลักษณะตัวอักษร โดยแบ่งออกได้ 4 กลุ่มดังต่อไปนี้ [8]

- Basic features เป็นคุณลักษณะพื้นฐานที่ได้จากแพคเกจข้อมูลที่มีการสื่อสารในเครือข่าย เช่น เวลาในการเชื่อมต่อ ชนิดของโปรโตคอล ชนิดของการให้บริการ และสถานะแพ็กเก็ต เป็นต้น มีทั้งหมด 9 คุณลักษณะ

- Content features เป็นคุณลักษณะที่เก็บรวบรวมข้อมูลที่แสดงให้เห็นถึงพฤติกรรมน่าสงสัย เช่น ความผิดพลาดในการถือกรอกอิน หรือการใช้คำสั่ง "su" เป็นต้น มีทั้งหมด 13 คุณลักษณะ

- Traffic features เป็นคุณลักษณะที่เก็บรวบรวมข้อมูลที่แสดงลักษณะของการสื่อสาร เช่น จำนวนในการครั้งในการเชื่อมต่อเข้าสู่ระบบเมื่อผ่านไประยะเวลา 2 วินาที เป็นต้น มีทั้งหมด 9 คุณลักษณะ

- Host based features เป็นคุณลักษณะที่เก็บรวบรวมข้อมูลที่แสดงลักษณะของการสื่อสารไปยังเครื่องปลายทางเครื่องเดิมตลอดเวลา เช่น จำนวนครั้งในการเชื่อมต่อไปยังเครื่องปลายทางเครื่องเดิม เป็นต้น มีทั้งหมด 10 คุณลักษณะ

ส่วนที่สองเป็นส่วนการนำ SOM มาประยุกต์ใช้ในการแยกประเภทและตรวจจับการบุกรุก แล้วส่งต่อไปยังส่วนสุดท้ายเพื่อดำเนินการแจ้งเตือนไปยังเครื่องคอมพิวเตอร์

3.2. การเตรียมข้อมูลอินพุตเวกเตอร์

ในการเตรียมข้อมูลอินพุตเวกเตอร์เพื่อใช้กับ SOM นั้นงานวิจัยนี้ได้กำหนดการทดลองบนพื้นฐานของประเภทการบุกรุกแบบ Denial of Service กับ Probing

และพฤติกรรมปกติ เนื่องจากเป็นกลุ่มข้อมูลที่มีจำนวนมากที่สุด มีทั้งสิ้น 492,843 เรคคอร์ด

และจากคุณลักษณะ (Feature) ทั้ง 41 คุณลักษณะ (Feature) ดังกล่าวนั้นคุณลักษณะ (Feature) ชนิดของโปรโตคอล (เช่น TCP) ชนิดการให้บริการ (เช่น http) และสถานะแพ็กเก็ต เป็นคุณลักษณะ (Feature) แบบสัญลักษณ์ (Symbolic) ซึ่งต้องทำการจัดให้อยู่ลักษณะของตัวเลขก่อน จากนั้นจึงแบ่งข้อมูลออกเป็น 2 กลุ่มเพื่อใช้ในการสอนระบบจำนวน 394,274 เรคคอร์ด และใช้ในการทดสอบระบบจำนวน 98,569 เรคคอร์ด ดังตารางที่ 1

Attack Type	Train	Test
normal.	77,822	19,456
neptune.	85,761	21,440
smurf.	224,632	56,158
back.	1,762	441
satan.	1,271	318
ipsweep.	998	249
portsweep.	832	208
teardrop.	783	196
pod.	211	53
land.	17	4
nmap.	185	46
Total	394,274	98,569

ตารางที่ 1 แสดงกลุ่มข้อมูลที่ใช้ทดลอง

3.3. การสร้างแผนภาพ SOM แบบมัดดีลำดับชั้น

การทดลองจะใช้ SOM ขนาด 30x30 โดยมีจำนวนโหนด 900 โหนดเพื่อใช้ในการสร้าง SOM ลำดับชั้นที่หนึ่ง ส่วนลำดับชั้นที่เพิ่มใช้ SOM ขนาด 5x5 โดยมีจำนวนโหนด 25 โหนด การเพิ่มลำดับชั้นพิจารณาจากอัตราส่วนของประเภทการบุกรุกที่ทับซ้อนกันในโหนดเดียวกัน ถ้าหากค่าอัตราส่วนที่มีค่ามากที่สุดที่คำนวณได้ และมีค่ามากกว่า β ที่กำหนด จึงเพิ่มลำดับชั้น และทำไปเรื่อยๆ จนกว่าครบทุกโหนดที่มีประเภทการบุกรุกที่ทับซ้อนกัน ดังสมการที่ 6 และ 7 โดยในงานวิจัยนี้ใช้ค่า $\beta = 0.9$ ผลการทดลองของงานวิจัยนี้ได้จากการเปรียบเทียบกันระหว่าง SOM ลำดับชั้นเดียว กับ SOM หลายลำดับชั้น

$$CR = \left(\frac{N_{Select}}{N_T} \right) \quad (6)$$

เมื่อ

N_{Select} แทนจำนวนประเภทการบุกรุกที่เลือกในโหนด

N_{Tn} แทนจำนวนประเภทการบุกรุกทั้งหมดในโหนด

สมการที่ 7 แสดงการหาค่าเงื่อนไขในการเพิ่มลำดับชั้นของแผนภาพ SOM

$$AddLayer(CR_{Max}) = \begin{cases} true, & \text{if } CR_{Max} < \beta \\ false, & \text{otherwise} \end{cases} \quad (7)$$

เมื่อ

CR_{Max} แทนอัตราส่วนประเภทของการบุกรุกที่ทับซ้อนกันในโหนดเดียวกันที่มีค่ามากที่สุด

β แทนค่ากำหนดเงื่อนไขในการเพิ่มจำนวนลำดับชั้น

อัลกอริทึมในการสร้าง SOM แบบหลายลำดับชั้นมีดังนี้

1. หาโหนดในแผนภาพชั้นที่หนึ่ง ที่มีข้อมูลที่ทับซ้อนกัน
2. คำนวณค่าอัตราส่วน CR ดังสมการที่ 6 ในแต่ละประเภทของการบุกรุกและเปรียบเทียบค่าที่สูงที่สุด
3. พิจารณาเพิ่มลำดับชั้น ดังสมการที่ 7
 - 3.1. ถ้าค่า $CR_{Max} \geq \beta$ ไม่เพิ่มลำดับชั้น
 - 3.2. ถ้าค่า $CR_{Max} < \beta$ เพิ่มลำดับชั้น
4. ทำซ้ำในข้อที่ 1 ถึง 3 ตามลำดับจนครบทุกโหนดที่มีข้อมูลที่ทับซ้อน

เมื่อดำเนินการตามอัลกอริทึมเสร็จ จะได้ SOM แบบหลายลำดับชั้นเฉพาะ โหนดที่มีข้อมูลที่ทับซ้อนเท่านั้น

3.4. การกำหนดค่า Detection Rate และ False Positive

สมการที่ 8 แสดงการหาค่าเปอร์เซ็นต์การตรวจจับหรือเปอร์เซ็นต์ Detection rate

$$\%Detected = \left(\frac{N_{attack} - N_{missed}}{N_{attack}} \right) * 100 \quad (8)$$

เมื่อ

N_{attack} แทนจำนวนของประเภทการบุกรุกทั้งหมดในข้อมูลทดสอบ

N_{missed} แทนจำนวนของการตรวจจับผิดพลาด

สมการที่ 9 ใช้สำหรับการหาค่าเปอร์เซ็นต์ False Positive โดยค่าเปอร์เซ็นต์ False Positive คือค่าเปอร์เซ็นต์

ของจำนวนของพฤติกรรมปกติ แต่ระบบตรวจจับแจ้งว่าเป็นพฤติกรรมการบุกรุก

$$\%FalsePositive = \left(\frac{N_{false}}{N_{normal}} \right) * 100 \quad (9)$$

เมื่อ

N_{false} แทนจำนวนทั้งหมดของ false positive ที่พบในการทดสอบ

N_{normal} แทนจำนวน normal ทั้งหมดในข้อมูลทดสอบ

4. ผลการทดลอง

ในตารางที่ 2 เป็นผลการทดลองคำนวณหาค่าเปอร์เซ็นต์ Detection rate ของ SOM แบบลำดับชั้นเดียวในการตรวจจับพฤติกรรมการบุกรุกจำนวน 11 ประเภท โดยรวมพฤติกรรมปกติด้วยแต่การคำนวณหาค่าเปอร์เซ็นต์ Detection rate ของ SOM แบบลำดับชั้นเดียวนั้น งานวิจัยนี้คำนวณเฉพาะค่าของโหนดที่ไม่มีข้อมูลการบุกรุกที่ทับซ้อนกัน ส่วนการบุกรุกที่ทับซ้อนกัน งานวิจัยนี้ถือว่าเป็นการตรวจจับที่ไม่แน่นอน ไม่สามารถระบุว่าเป็นพฤติกรรมการบุกรุกใดที่แน่นอนได้ ซึ่งงานวิจัยนี้ได้เสนอการแก้ปัญหาดังกล่าวนี้โดยเพิ่มลำดับชั้นของ SOM ดังผลการทดลองในตารางที่ 3

Attack Type	Detection Rate
Normal	99.72%
Neptune	99.97%
Smurf	99.74%
Back	99.50%
Satan	93.79%
Ipsweep	93.10%
Portsweep	66.67%
Teardrop	98.96%
Pod	100.00%
Land	100.00%
Nmap	50.00%

ตารางที่ 2 Detection Rate ของ SOM แบบลำดับชั้นเดียว

ตารางที่ 3 เป็นผลการทดลองการคำนวณค่าเปอร์เซ็นต์ Detection rate ของ SOM แบบหลายลำดับชั้นซึ่งแก้ปัญหาข้อมูลทับซ้อนโดยการเพิ่มลำดับชั้นเฉพาะ

โหนดที่มีข้อมูลที่ทับซ้อน และกำหนดจำนวนของประเภท การบุกรุกทั้งหมดในข้อมูลทดสอบ (N_{attack}) ให้มีค่า เท่ากับจำนวนการบุกรุกทั้งหมดที่อยู่ในโหนดที่ทำการเพิ่ม ลำดับชั้นนั้น

ตารางที่ 4 เป็นการคิดคำนวณค่าเปอร์เซ็นต์ False positive ซึ่งได้จากการตรวจจับพฤติกรรมที่ผิดพลาดของ ระบบที่ตรวจจับพฤติกรรมปกติ แต่รายงานว่าเป็น พฤติกรรมการบุกรุก โดยในส่วนของ SOM แบบลำดับชั้น เดียวนั้น งานวิจัยนี้คำนวณโดยคิดเฉพาะโหนดที่มีข้อมูล ไม่ทับซ้อนกัน ส่วนข้อมูลที่ทับซ้อนคิดเป็นเปอร์เซ็นต์การ Error rate และส่วน SOM แบบหลายลำดับชั้นจะคิดค่า เปอร์เซ็นต์ False positive รวมกับลำดับชั้นที่เพิ่มด้วย

Attack Type	Detection Rate
Normal	99.94%
Neptune	100.00%
Smurf	99.86%
Back	100.00%
Satan	96.92%
Ipsweep	94.27%
Portsweep	92.67%
Teardrop	100.00%
Pod	100.00%
Land	100.00%
Nmap	97.46%

ตารางที่ 3 Detection Rate ของ SOM แบบหลายลำดับชั้น

SOM	False positive
ลำดับชั้นเดียว	0.23%
หลายลำดับชั้น	0.16%

ตารางที่ 4 False positive แผนภาพ SOM

5. บทสรุป

จากการทดลองโดยการประยุกต์ใช้ SOM มาทำการ แยกประเภทการบุกรุกนั้น ผลปรากฏว่าในการทดลองโดย ใช้ SOM แบบลำดับชั้นเดียวนั้น มีข้อมูลที่ทับซ้อนกันใน โหนดเดียวกัน จึงทำให้ประสิทธิภาพในการตรวจจับการ บุกรุกระบบเครือข่าย หรือ %Detection rate จึงมีค่าต่ำ โดยเฉพาะการบุกรุกประเภท Nmap จากตารางที่ 2 และ

ค่า %false positive มีค่าสูงจากตารางที่ 4 ทำให้ระบบ ตรวจจับการบุกรุกมีความน่าเชื่อถือต่ำ

เมื่อได้ทำการประยุกต์ SOM แบบหลายลำดับชั้นมาใช้ ในการเพิ่มลำดับชั้นเฉพาะ โหนดที่มีข้อมูลที่ทับซ้อนกัน ผลปรากฏว่า SOM แบบหลายลำดับชั้นมีประสิทธิภาพใน การตรวจจับการบุกรุก หรือ %Detection rate มีค่าสูงขึ้น เฉลี่ยประมาณ 7.24 เปอร์เซ็นต์ และค่า %false positive มี ค่าต่ำลงประมาณ 15.77 เปอร์เซ็นต์ ซึ่งส่งผลทำให้ระบบ ตรวจจับการบุกรุกมีความน่าเชื่อถือมากยิ่งขึ้น ต่อไปน่าจะ ดำเนินการวิจัยในส่วนของการเลือกกลุ่มของคุณลักษณะ (Feature) และเมื่อนำไปใช้ในการตรวจจับพฤติกรรม การบุกรุกประเภทต่างๆ แล้ว ประสิทธิภาพในการตรวจจับการ บุกรุกนั้นควรมีประสิทธิภาพเป็นอย่างดี

6. เอกสารอ้างอิง

- [1] P. Dokas, L. Ertöz, V. Kumar, A. Lazarevi, J. Srivastava and P. Tan, "Data mining for network intrusion detection," In Proceedings NSF Workshop on Next Generation Data Mining, p.21-30, 2002.
- [2] W. Lee and S.J. Stolfo, "Data mining approaches for intrusion detection," In Proceedings of the 7th USENIX Security Symposium, San Antonio, TX, 1998.
- [3] S. Lee and D. Heinbuch, "Training a neural-network based intrusion detector to recognize novel attacks," IEEE Transactions on Systems, Man & Cybernetics, Part A (Systems & Humans), 31(4):294-9, July 2001.
- [4] B. C. Rhodes, J. A. Mhahffey and J. D. Cannady, "Multiple self-organizing maps for intrusion detection," In Proceedings of the 23rd National Information Systems Security Conference, 2000.
- [5] H. Kayacik, A. Zincir-Heywood and M. Heywood, "On the capability of an SOM based Intrusion Detection System," In Proceedings. IEEE Int. Joint Conf. Neural Networks (IJCNN'03), pp. 1808-1813.
- [6] Suseela T., Qiuming A. Zhu and Julie Huff, "Hierarchical Kohonen Net for Anomaly Detection in Network Security," In IEEE Transactions on Systems, Vol.35, No.2, 2005.
- [7] T. Kohonen, "Self-Organizing Maps," 3rd edition, Springer Springer-Verlag, 2001
- [8] S. Stolfo et al., The Third International Knowledge Discovery and Data Mining Tools Competition <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>

ประวัติผู้เขียน

ชื่อ-สกุล	นายสุรพล โรจนประดิษฐ์
วันเดือนปีเกิด	วันที่ 6 มิถุนายน พ.ศ. 2518 ณ จังหวัดสุพรรณบุรี
ที่อยู่	450 ม.6 ถนน สุพรรณ-ชัยนาท ตำบล ย่านยาว อำเภอสามชุก จังหวัด สุพรรณบุรี 72130
ประวัติการศึกษา	
พ.ศ. 2541	สำเร็จการศึกษาวิศวกรรมศาสตรบัณฑิต สาขาวิศวกรรมคอมพิวเตอร์ ศูนย์กลางสถาบันเทคโนโลยีราชมงคล
พ.ศ. 2545	เข้าศึกษาต่อในระดับวิศวกรรมศาสตรมหาบัณฑิต สาขาวิศวกรรมคอมพิวเตอร์ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
ประสบการณ์ทำงาน	
พ.ศ. 2542-ปัจจุบัน	เข้ารับราชการตำแหน่งอาจารย์ประจำแผนกวิชาเทคนิคคอมพิวเตอร์ คณะวิชาไฟฟ้า สถาบันเทคโนโลยีราชมงคล วิทยาเขตสุพรรณบุรี