

สำนักหอสมุดกลาง พระจอมเกล้าลาดกระบัง

ระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่

NETWORK INTRUSION DETECTION SYSTEM USING ROUGH FUZZY



พ.
๕๓๒๒
๒๕๕๐

เลขหมู่.....
เลขทะเบียน.....**74494**
วัน,เดือน,ปี.....- 2 ต.ค. 2550

b.....11๙.๒๕๓๓
i.....

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิศวกรรมศาสตรมหาบัณฑิต

สาขาวิชาวิศวกรรมคอมพิวเตอร์

บัณฑิตวิทยาลัย

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

พ.ศ.๒๕๕๐

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

NETWORK INTRUSION DETECTION SYSTEM USING ROUGH FUZZY



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

2007

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



COPYRIGHT 2007

SCHOOL OF GRADUATE STUDIES

KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

หัวข้อวิทยานิพนธ์	ระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิชซี่
นักศึกษา	นายธรรมกร ครองไตรภพ
รหัสนักศึกษา	45061220
ปริญญา	วิศวกรรมศาสตรมหาบัณฑิต
สาขาวิชา	วิศวกรรมคอมพิวเตอร์
พ.ศ.	2549
อาจารย์ที่ปรึกษาวิทยานิพนธ์	รศ.ดร. เอียน ปิ่นเงิน

บทคัดย่อ

ปัจจุบันระบบตรวจจับการบุกรุกเครือข่าย โดยใช้เทคนิคการเรียนรู้ของเครื่องจักร (machine learning) เป็นหัวข้องานวิจัยที่แพร่หลาย เนื่องจากสามารถตรวจจับได้ทั้งการบุกรุกแบบ misuse และ anomaly งานวิจัยนี้ได้เสนอระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิชซี่ ระบบที่นำเสนอการใช้ฟิชซี่เซตและราฟเซตในการสร้างกฎที่ใช้จำแนกการบุกรุกสำหรับจำแนกพฤติกรรมบุกรุกและพฤติกรรมปกติ โดยแบ่งการทำงานออกเป็นสี่ขั้นตอนหลัก ขั้นตอนแรกเป็นการสร้างระบบตัดสินใจโดยใช้ข้อมูลจากฐานข้อมูลของ KDD ขั้นตอนที่สองแบ่งช่วงลักษณะเด่นที่มีความต่อเนื่องของข้อมูลโดยใช้ฟิชซี่เซต ขั้นตอนที่สามนำข้อมูลที่ได้มากำจัดคุณลักษณะที่ไม่มีความจำเป็นในการจำแนกและสร้างกฎที่ใช้จำแนกการบุกรุกสำหรับจำแนกพฤติกรรมบุกรุกและพฤติกรรมปกติโดยใช้ราฟเซต และขั้นตอนสุดท้ายทดสอบประสิทธิภาพของกฎที่ใช้จำแนกการบุกรุก จากการทดลองระบบด้วยฐานข้อมูล KDD 99 พบว่า Detection rate เฉลี่ยเท่ากับ 99.37% และ False negative rate เท่ากับ 0.079% ซึ่งมีความแม่นยำสูงกว่าระบบที่ใช้นิวรอลเน็ตเวิร์คและคิซิทันตรี

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Thesis Title	Network Intrusion Detection System Using Rough Fuzzy
Student	Mr. Thammakorn Krongtripop
Student ID.	45061220
Degree	Master of Engineering
Program	Computer Engineering
Year	2006
Thesis Advisor	Assoc. Prof. Dr. Ouen Pinngern

ABSTRACT

Recently machine learning-based intrusion detection approaches have been subjected to extensive researches because they can detect both misuse and anomaly. In this research, we present an intrusion detection system based on a rough-fuzzy approach. For the proposed system, we use fuzzy set and rough set to create a set of intrusion rules for classifying normal and abnormal behaviors. The process consists of four stages. First, create the decision system using data from KDD database. Second, discretize the continuous data to subintervals using fuzzy set. Third, reduce the redundancy, inconsistency of initial information table and generate intrusion rules using rough set. Finally, test for efficiency of intrusion rules using KDD 99 data set. From the experiments, the average of detection rate is 99.37% and of false negative rate is 0.079% which is higher accuracy than those systems using neural network and decision tree.

กิตติกรรมประกาศ

วิทยานิพนธ์นี้สำเร็จลุล่วงด้วยดีเนื่องจากกำลังใจและพระคุณอันหาที่สุดมิได้ จากคุณพ่อคุณแม่ และพี่สาวผู้ล่วงลับ ข้าพเจ้าขอสำนึกในพระคุณนี้อย่างเป็นที่สุด

วิทยานิพนธ์นี้จะไม่สำเร็จลุล่วงหากปราศจากแรงผลักดัน และคำแนะนำที่มีประโยชน์ของ รศ.ดร. เอื้อน ปิ่นเงิน ผู้ควบคุมวิทยานิพนธ์ ข้าพเจ้าขอกราบขอบพระคุณเป็นอย่างสูง

ข้าพเจ้าขอกราบเท้า คุณครูและอาจารย์ทุกท่านตั้งแต่เล็กจนเติบโตใหญ่ ที่ได้มอบวิชาความรู้ให้แก่ข้าพเจ้า รวมทั้งคำสั่งสอนและอบรมให้ข้าพเจ้าเป็นคนดี ข้าพเจ้าขอกราบขอบพระคุณเป็นอย่างสูง

ข้าพเจ้าขอขอบพระคุณ ภาควิชาวิศวกรรมคอมพิวเตอร์ และสำนักวิจัยการสื่อสารและเทคโนโลยีสารสนเทศ (ReCCIT) ที่ได้สนับสนุนเครื่องมือ ตลอดจนข้อมูล และหนังสือต่างๆ ที่ใช้ในการทำวิจัย ข้าพเจ้าขอกราบขอบพระคุณเป็นอย่างสูง

ข้าพเจ้าขอขอบคุณสำหรับกำลังใจ คำแนะนำ และประสบการณ์ที่ดีจากพี่ ๆ และเพื่อน ๆ นักศึกษา ป.โททุกท่าน หัวปีที่ลาดกระบังข้าพเจ้าจะไม่มีวันลืม และขอขอบคุณนายพรเทพ โรจนวสุ และนายไพฑูรย์ ศรีนิล ที่ช่วยแก้ไขภาษาในการพิมพ์บทความตีพิมพ์ ข้าพเจ้าขอขอบคุณ

สุดท้ายนี้คุณค่าและประโยชน์อันพึงมีจากวิทยานิพนธ์ฉบับนี้ ข้าพเจ้าขอมอบให้กับผู้มีพระคุณทุกท่าน หากวิทยานิพนธ์ฉบับนี้มีข้อผิดพลาดประการใดข้าพเจ้าขอน้อมรับไว้เพียงผู้เดียว

ธรรมกร ครองไตรภพ

สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	I
บทคัดย่อภาษาอังกฤษ.....	II
กิตติกรรมประกาศ.....	III
สารบัญ.....	IV
สารบัญตาราง.....	VII
สารบัญภาพ.....	IX
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 เนื้อหาวิทยานิพนธ์.....	3
บทที่ 2 ระบบตรวจจับการบุกรุกและงานวิจัยที่เกี่ยวข้อง.....	4
2.1 ระบบตรวจจับการบุกรุก.....	4
2.1.1 ความหมายของระบบตรวจจับการบุกรุก.....	4
2.1.2 ประเภทของระบบตรวจจับการบุกรุก.....	6
2.1.3 การทำงานของระบบตรวจจับการบุกรุก.....	8
2.1.4 ความสำคัญของระบบตรวจจับการบุกรุก.....	14
2.2 งานวิจัยที่เกี่ยวข้อง.....	15
2.3 สรุป.....	16
บทที่ 3 ฟิชชั่นเช็ต.....	17
3.1 ทฤษฎีฟิชชั่นเช็ต.....	17
3.2 การแทนข้อมูลในฟิชชั่นเช็ต.....	18
3.3 ฟังก์ชันความเป็นสมาชิก.....	19
3.4 สรุป.....	23

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ(ต่อ)

	หน้า
บทที่ 4 ราวเฟ้ต.....	24
4.1 การแสดงค่าความรู้.....	24
4.1.1 ระบบสารสนเทศ.....	25
4.1.2 ระบบการตัดสินใจ.....	25
4.2 ความคล้ายกันของวัตถุ.....	26
4.3 การประมาณค่าของเซ้ต.....	27
4.4 เซ้ตที่มีคุณลักษณะน้อยที่สุดที่ยังคงสามารถจำแนกวัตถุได้.....	29
4.5 การขึ้นต่อกันของคุณลักษณะ.....	32
4.6 กฎการตัดสินใจ.....	33
4.7 สรุป.....	34
บทที่ 5 การออกแบบระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราวเฟ้ต.....	35
5.1 ขั้นตอนการทำงานของระบบ.....	35
5.2 ข้อมูล KDD Cup 1999.....	35
5.2.1 ลักษณะข้อมูลของ KDD Cup 1999.....	36
5.2.2 การเตรียมข้อมูลอินพุต.....	39
5.3 การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราวเฟ้ต.....	39
5.3.1 กระบวนการสร้างระบบการตัดสินใจ.....	40
5.3.2 กระบวนการแปลงข้อมูลเบื้องต้น.....	42
5.3.3 กระบวนการสร้างกฎที่ใช้จำแนกการบุกรุก.....	44
5.3.4 กระบวนการทดสอบกฎที่ใช้จำแนกการบุกรุก.....	47
5.4 สรุป.....	47
บทที่ 6 ผลการทดลอง.....	49
6.1 การเตรียมเครื่องมือและข้อมูลที่ใช้ในการทดลอง.....	49
6.1.1 เครื่องมือที่ใช้ในการทดลอง.....	49
6.1.2 การเตรียมข้อมูลที่ใช้ในการทดลอง.....	49
6.2 การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราวเฟ้ต.....	53
6.3 การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิวรอลเน็ตเวิร์ค.....	59

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ(ต่อ)

	หน้า
6.4 การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ดิจิทัลรี.....	65
6.5 เปรียบเทียบผลการตรวจจับการบุกรุกผ่านเครือข่ายของแต่ละระบบ.....	71
บทที่ 7 สรุปผลการวิจัยและข้อเสนอแนะ.....	73
7.1 สรุปผลการวิจัย.....	73
7.2 ข้อเสนอแนะ.....	73
บรรณานุกรม.....	74
ภาคผนวก งานวิจัยที่ได้รับการตีพิมพ์.....	76
ประวัติผู้เขียน.....	84

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญตาราง

ตารางที่	หน้า
4.1 รูปแบบของระบบสารสนเทศและระบบการตัดสินใจ.....	24
4.2 ตัวอย่างของระบบสารสนเทศของคนใช้.....	25
4.3 ตัวอย่างของระบบการตัดสินใจ.....	26
4.4 ตัวอย่างของระบบตัดสินใจที่ยังไม่ได้ลดคุณลักษณะที่มีชื่อว่า <i>Hiring</i>	30
4.5 Discernibility matrix ของระบบตัดสินใจ <i>Hiring</i>	31
5.1 คุณลักษณะพื้นฐาน.....	36
5.2 Content features.....	37
5.3 Traffic features.....	37
5.4 Host based features.....	38
5.5 แสดงกลุ่มข้อมูลที่ใช้ทดลอง.....	39
5.6 แสดงตัวอย่างข้อมูล.....	41
5.7 แสดงตัวอย่างข้อมูลหลังผ่านการแปลงข้อมูลเบื้องต้น.....	43
5.8 แสดงตัวอย่าง reducts.....	44
5.9 แสดงตัวอย่างกฎที่ใช้จำแนกการบุกรุก.....	46
6.1 แสดงกลุ่มข้อมูลที่ใช้ทดลอง.....	50
6.2 แสดงตัวอย่างข้อมูลที่ใช้ในการสอนและทดสอบระบบ.....	50
6.3 แสดงรูปแบบการแบ่งข้อมูลที่ใช้ทดลอง.....	51
6.4 แสดงการกำหนดค่าตัวเลขแทน Protocol feature.....	52
6.5 แสดงการกำหนดค่าตัวเลขแทน Service feature.....	52
6.6 แสดงการกำหนดค่าตัวเลขแทน Flag feature.....	53
6.7 แสดงตัวอย่างข้อมูลที่ถูกแปลงค่าคุณลักษณะที่มีความต่อเนื่องให้มีค่าอยู่ระหว่าง 0 ถึง 1.....	54
6.8 แสดงตัวอย่างข้อมูลที่ถูกแปลงค่าคุณลักษณะที่มีความต่อเนื่องให้เป็นช่วงโดยใช้ fuzzy spaces.....	55
6.9 แสดงจำนวนคุณลักษณะ ค่าการขึ้นต่อกันของคุณลักษณะและกฎที่ใช้จำแนกพฤติกรรม การบุกรุกที่ได้จากการเรียนรู้ของระบบ.....	56
6.10 แสดงตัวอย่างกฎที่ใช้จำแนกการบุกรุกที่ได้จากการเรียนรู้ของระบบ.....	57
6.11 Detection rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิชซี่.....	58
6.12 False negative rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิชซี่.....	58
6.13 แสดงตัวอย่างข้อมูลที่ถูกแปลงค่าคุณลักษณะที่มีความต่อเนื่องให้มีค่าอยู่ระหว่าง 0 ถึง 1.....	60

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญตาราง(ต่อ)

ตารางที่	หน้า
6.14 แสดงตัวอย่างข้อมูลหลังแปลงคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์ให้ตัวเลข.....	61
6.15 พารามิเตอร์ในการสอนระบบเริ่มต้น.....	63
6.16 แสดงค่า Mean Square Error ที่ทดลองได้จากข้อมูลแต่ละชุด.....	63
6.17 Detection rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิวรอลเน็ตเวิร์ค.....	64
6.18 False negative rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิวรอลเน็ตเวิร์ค.....	64
6.19 แสดงตัวอย่างข้อมูลที่ถูกแปลงค่าคุณลักษณะที่มีความต่อเนื่องให้มีค่าอยู่ระหว่าง 0 ถึง 1.....	66
6.20 แสดงตัวอย่างข้อมูลหลังแปลงคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์ให้ตัวเลข.....	67
6.21 พารามิเตอร์ในการสร้างคีชีซันทรี.....	68
6.22 แสดงจำนวน โหนด คุณลักษณะและกฎที่ใช้จำแนกพฤติกรรมการบุกรุกที่ได้จากการทดลอง.....	69
6.23 แสดงตัวอย่างกฎที่ใช้จำแนกการบุกรุกที่ได้จากการเรียนรู้ของระบบ.....	69
6.24 Detection rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้คีชีซันทรี.....	69
6.25 False negative rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้คีชีซันทรี.....	70

สารบัญรูป

รูปที่	หน้า
2.1 ระบบการตรวจจับการบุกรุก.....	5
2.2 ขอบเขตที่เหลื่อมกันของ IDS กับระบบทำให้เกิด False positive และ False negative.....	6
2.3 Anomaly Detection Model.....	7
2.4 Misuse Detection Model.....	7
2.5 การทำงานของระบบตรวจจับการบุกรุก.....	8
3.1 กราฟแสดงค่าความเป็นสมาชิกของ x ที่มีต่อเซต A กรณีที่เป็นคริสป์เซต (crisp set).....	18
3.2 กราฟแสดงค่าความเป็นสมาชิกของเซต A “ความสูง” กรณีที่เป็นฟัซซี่เซต.....	19
3.3 กราฟแสดงฟังก์ชันความเป็นสมาชิกของตัวแปร “ความสูง” ที่ประกอบด้วย 3 เทอมเซต.....	20
3.4 กราฟแสดงฟังก์ชันความเป็นสมาชิกแบบรูปสามเหลี่ยม.....	20
3.5 กราฟแสดงฟังก์ชันความเป็นสมาชิกแบบรูปสี่เหลี่ยมคางหมู.....	21
3.6 กราฟแสดงฟังก์ชันความเป็นสมาชิกแบบรูปตัว S.....	22
3.7 กราฟแสดงฟังก์ชันความเป็นสมาชิกแบบรูปประมงคว่ำ.....	22
3.8 กราฟแสดงฟังก์ชันความเป็นสมาชิกแบบรูป Gaussian.....	23
4.1 การประมาณค่าเซตคนไข้ที่เดินได้ โดยใช้ Age และ LEMS เป็นคุณลักษณะที่ใช้ในการสร้างเงื่อนไข กลุ่มของคนไข้ที่มีอาการคล้ายกันจะอยู่ใน regions เดียวกัน.....	29
5.1 แสดงขั้นตอนการทำงานของระบบ.....	35
5.2 การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้กราฟฟัซซี่.....	40
5.3 Fuzzy spaces ของแต่ละคุณลักษณะ.....	42
6.1 โครงสร้าง BP network ของระบบตรวจจับการบุกรุก.....	59
6.2 กราฟแสดงค่า Detection rate เหนือของระบบตรวจจับการบุกรุกต่างๆ.....	71
6.3 กราฟแสดงค่า False negative rate เหนือของระบบตรวจจับการบุกรุกต่างๆ.....	71

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

ปัจจุบันการตรวจจับการบุกรุกเป็นสิ่งที่ยากเป็นอย่างมาก เนื่องจากมีการบุกรุกระบบคอมพิวเตอร์เพิ่มมากขึ้นตามจำนวนคอมพิวเตอร์ที่ต่อเข้าสู่เครือข่าย ก่อให้เกิดความเสียหายกับองค์กรทางธุรกิจต่างๆ ดังนั้นการตรวจจับการบุกรุกเครือข่ายจึงกลายเป็นหัวข้อการวิจัยที่สำคัญ เนื่องจากไม่มีระบบตรวจจับการบุกรุกใดที่ไม่มีช่องโหว่ และความสนใจในการพัฒนาระบบตรวจจับการบุกรุกเครือข่ายได้มุ่งประเด็นไปที่การใช้เทคนิคการเรียนรู้ของเครื่องจักร (machine learning) ซึ่งเป็นหัวข้องานวิจัยที่แพร่หลาย เนื่องจากสามารถตรวจจับได้ทั้งการบุกรุกแบบ misuse และ anomaly โดยเทคนิคที่ถูกนำมาใช้อาทิเช่น นิวรอลเน็ตเวิร์คที่มีความถูกต้องในการตรวจจับสูง แต่มีปัญหาที่กฎที่ใช้จำแนกการบุกรุก (intrusion rules) ไม่สามารถอธิบายได้ และคิซึซันทรีที่มีกฎที่ใช้จำแนกการบุกรุกสามารถอธิบายได้ แต่มีปัญหาที่มีความถูกต้องในการตรวจจับต่ำ

งานวิจัยนี้นำเสนอวิธีการตรวจจับการบุกรุกผ่านเครือข่ายแบบ misuse โดยใช้ราฟฟิซึซึ เนื่องจากระบบตรวจจับการบุกรุกส่วนมากถูกสร้างโดยใช้รูปแบบพฤติกรรมการบุกรุกที่ได้รับรองและรายงานโดยผู้เชี่ยวชาญ และรูปแบบพฤติกรรมที่มีอยู่ในปัจจุบันเพียงพอที่จะใช้ตรวจจับการบุกรุกใหม่ๆ ได้ นอกจากนี้วิธีการตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซึซึยังสามารถสร้างกฎที่ใช้จำแนกการบุกรุกที่สามารถอธิบายได้และมีความถูกต้องในการตรวจจับสูงซึ่งประกอบด้วย 4 ขั้นตอน โดยขั้นตอนแรกจะทำการสร้างระบบตัดสินใจโดยใช้ข้อมูลจากฐานข้อมูลของ KDD Cup 1999 ขั้นตอนที่สองแบ่งช่วงลักษณะเด่นที่มีความต่อเนื่องของข้อมูลโดยใช้ฟิซึซึเซต ขั้นตอนที่สามกำจัดคุณลักษณะที่ไม่มีความจำเป็นในการจำแนกและสร้างกฎที่ใช้จำแนกการบุกรุกสำหรับจำแนกพฤติกรรมบุกรุกและพฤติกรรมปกติโดยใช้ราฟฟิซึเซต และขั้นตอนสุดท้ายทดสอบประสิทธิภาพของกฎที่ใช้จำแนกการบุกรุก โดยทำการทดลองบนพื้นฐานการบุกรุกแบบเครือข่ายแบบ Denial of Service (DoS) Probing และพฤติกรรมปกติ ด้วยวิธีการที่นำเสนอในวิทยานิพนธ์ทำให้ประสิทธิภาพการตรวจจับผู้บุกรุกดีกว่าวิธีนิวรอลเน็ตเวิร์ค (neural network) และคิซึซันทรี (decision tree) [13]

1.2 จุดมุ่งหมายและวัตถุประสงค์ของการศึกษา

เอกสารนี้เป็นเพื่อศึกษาวิธีการลดแอตทริบิวหรือคุณลักษณะที่ไม่จำเป็นต่อการจำแนกพฤติกรรมบุกรุกไม่ว่ากรณีใดและพฤติกรรมปกติโดยใช้ราฟฟิซึซึ และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2. เพื่อศึกษาแนวทางในการพัฒนาระบบตรวจจับการบุกรุกและนำเสนอวิธีใหม่ที่สามารถวิจัยและพัฒนาขึ้น
3. เพื่อศึกษาการจำแนกพฤติกรรมบุกรุกแบบ DoS Probing และพฤติกรรมปกติโดยใช้วิธีกราฟฟิชชี
4. เพื่อเปรียบเทียบประสิทธิภาพการจำแนกพฤติกรรมบุกรุกแบบ DoS Probing และพฤติกรรมปกติระหว่างวิธีนิรอลเน็ตเวิร์คและดัชนีชั้นตรีกับวิธีใหม่ที่สามารถวิจัยนำเสนอ

1.3 แนวคิดที่ใช้ในงานวิจัย

ระบบตรวจจับการบุกรุกที่ดีจะต้องมีกฎที่ใช้จำแนกการบุกรุกที่สามารถอธิบายได้และมีความถูกต้องในการตรวจจับ (detection rate) สูง แต่ระบบตรวจจับการบุกรุกที่มีอยู่ไม่สามารถสร้างกฎที่ใช้จำแนกการบุกรุกที่สามารถทำทั้งสองอย่างพร้อมกันได้ [16] เช่น ระบบตรวจจับการบุกรุกที่ใช้นิรอลเน็ตเวิร์คสามารถสร้างกฎที่ใช้จำแนกการบุกรุกที่มีความถูกต้องในการตรวจจับสูง แต่ไม่สามารถอธิบายได้ ส่วนระบบตรวจจับการบุกรุกที่ใช้ดัชนีชั้นตรีสามารถสร้างกฎที่ใช้จำแนกการบุกรุกที่สามารถอธิบายได้ แต่มีความถูกต้องในการตรวจจับต่ำ

งานวิจัยนี้นำเสนอวิธีการตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้กราฟฟิชชี ซึ่งสามารถสร้างกฎที่ใช้จำแนกพฤติกรรมการบุกรุกที่สามารถอธิบายได้และมีความถูกต้องในการตรวจจับสูง ประกอบด้วย 4 ขั้นตอน โดยขั้นตอนแรกจะทำการสร้างระบบตัดสินใจโดยใช้ข้อมูลจากฐานข้อมูลของ KDD Cup 1999 ขั้นตอนที่สองแบ่งช่วงลักษณะเด่นที่มีความต่อเนื่องของข้อมูลโดยใช้ฟิชชีเซต ขั้นตอนที่สามกำจัดคุณลักษณะที่ไม่มีความจำเป็นในการจำแนกการบุกรุกและสร้างกฎที่ใช้จำแนกการบุกรุกสำหรับจำแนกพฤติกรรมการบุกรุกแบบ DoS Probing และพฤติกรรมปกติโดยใช้กราฟฟิชชี และขั้นตอนสุดท้ายทดสอบประสิทธิภาพของกฎที่ใช้จำแนกการบุกรุก ผลการทดลองตรวจจับการบุกรุกของวิธีการตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้กราฟฟิชชีเปรียบเทียบกับวิธีการตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิรอลเน็ตเวิร์คและดัชนีชั้นตรี [13]

1.4 การเปรียบเทียบระหว่างวิธีการที่นำเสนอกับวิธีที่มีอยู่เดิม

วิธีการเปรียบเทียบประสิทธิภาพของการตรวจจับการบุกรุกระบบเครือข่ายโดยใช้ค่าความถูกต้องในการตรวจจับ (detection rate) และการวัดค่าความผิดพลาดในการตรวจจับ (false negative rate) เมื่อเทียบกับหลักการที่มีอยู่เดิมแล้วในส่วน of ค่าความถูกต้องในการตรวจจับ จะให้ค่าความถูกต้องที่สูงกว่า และจะให้ค่าความผิดพลาดในการตรวจจับที่ต่ำกว่า

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1.5 ขอบเขตการวิจัย

1. ศึกษาการจำแนกพฤติกรรมการบุกรุกและพฤติกรรมปกติของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่
2. ข้อมูลพฤติกรรมการบุกรุกที่ใช้ในการทดสอบเป็นข้อมูลที่ได้จากการจำลองพฤติกรรมการบุกรุกบนระบบเครือข่าย โดยเลือกเฉพาะพฤติกรรมการบุกรุกประเภท Denial of Service กับ Probing และพฤติกรรมปกติ จากข้อมูลจำลองพฤติกรรมการบุกรุกชุดมาตรฐาน (KDD Cup 1999)

1.6 เนื้อหาของวิทยานิพนธ์

วิทยานิพนธ์ฉบับนี้ได้แบ่งเนื้อหาออกเป็น 7 บทคือ

บทที่ 1 กล่าวถึงความเป็นมาของงานวิจัย ความมุ่งหมายและวัตถุประสงค์ ขอบเขตของการวิจัย และขั้นตอนการศึกษา

บทที่ 2 กล่าวถึงระบบตรวจจับผู้บุกรุก ความหมายของระบบตรวจจับผู้บุกรุก ประเภทของระบบ การทำงานของระบบ และความสำคัญของระบบตรวจจับผู้บุกรุก

บทที่ 3 กล่าวถึงทฤษฎีของฟิซซี่เซต การแทนข้อมูลในฟิซซี่เซต และฟังก์ชันความเป็นสมาชิก

บทที่ 4 กล่าวถึงทฤษฎีของราฟฟิซซี่ การแสดงค่าความรู้ ความคล้ายกันของวัตถุ การประมาณค่าของเซต ซับเซตที่มีคุณลักษณะน้อยที่สุดที่ยังคงสามารถจำแนกวัตถุได้ การขึ้นต่อกันของคุณลักษณะ และกฎการตัดสินใจ

บทที่ 5 กล่าวถึงหลักการและวิธีการดำเนินงานวิจัย โดยจะเน้นถึงการเตรียมระบบเพื่อทำการทดสอบ

บทที่ 6 เป็นการทดลอง ผลการทดลอง

บทที่ 7 เป็นบทสรุปผลการวิจัยและข้อเสนอแนะ

บทที่ 2

ระบบตรวจจับการบุกรุกและงานวิจัยที่เกี่ยวข้อง

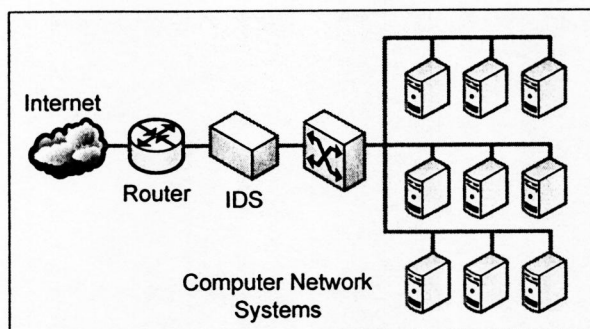
ในบทนี้จะอธิบายพื้นฐานของระบบตรวจจับการบุกรุกและงานวิจัยที่เกี่ยวข้อง โดยเนื้อหาประกอบด้วยระบบตรวจจับผู้บุกรุก ความหมายของระบบตรวจจับผู้บุกรุก ประเภทของระบบการทำงานของระบบตรวจจับผู้บุกรุก ความสำคัญของระบบตรวจจับผู้บุกรุก และงานวิจัยที่เกี่ยวข้องกับระบบตรวจจับการบุกรุก

2.1 ระบบตรวจจับการบุกรุก

ปัจจุบันมีการใช้งานระบบคอมพิวเตอร์อย่างแพร่หลาย หน่วยงานต่างๆ ได้นำระบบคอมพิวเตอร์มาใช้เพื่อพัฒนาศักยภาพการทำงานของตนให้มากขึ้น การทำงานของระบบคอมพิวเตอร์จะทำงานอย่างต่อเนื่อง และมีการออนไลน์ให้คนอื่นๆ เข้ามาใช้งานระบบบางส่วนด้วย ซึ่งปัญหาที่เกิดขึ้นตามมาคือปัญหาการบุกรุกเข้ามาสร้างความเสียหายให้กับระบบและอาจเข้ามาเพื่อขโมยข้อมูลที่สำคัญไป ปัญหาดังกล่าวจะเป็นปัญหากับผู้ดูแลระบบอย่างมากเนื่องจากผู้ดูแลระบบต้องคอยป้องกันและแก้ไขปัญหาต่างๆ อยู่เสมอๆ การทำงานของผู้ดูแลระบบนี้จะหนักมากหรือน้อยก็ขึ้นอยู่กับขนาดของระบบว่ามีขนาดใหญ่ และซับซ้อนมากน้อยเพียงใด ยิ่งระบบที่มีความซับซ้อนสูง มีความยุ่งยากในการดูแลมาก มีการออนไลน์ใช้งานอยู่ตลอดเวลา ผู้ดูแลระบบต้องดูแลระบบตลอดเวลาซึ่งในทางปฏิบัติแล้วเป็นไปได้ไม่ได้นั่นเอง อีกทั้งปัญหาบางปัญหาผู้ดูแลระบบไม่สามารถตรวจสอบด้วยตัวเองได้เลย ปัญหาการดูแลระบบจึงซับซ้อนขึ้นทุกวัน ทางออกที่จะช่วยแก้ปัญหาเหล่านี้ คือการมีผู้ช่วยมาช่วยดูแลระบบตลอดเวลาจะทำให้สามารถตรวจพบสิ่งที่ผู้ดูแลระบบไม่สามารถตรวจพบได้ คอยแจ้งเตือนให้กับผู้ดูแลระบบยามมีเหตุการณ์ผิดปกติเกิดขึ้น และในบางครั้งก็อาจสามารถแก้ไขปัญหาเบื้องต้น ได้ด้วย ผู้ช่วยที่สามารถช่วยแบ่งเบาภาระของผู้ดูแลระบบได้อย่างมากก็คือ ระบบตรวจจับการบุกรุก (Intrusion Detection System)

2.1.1 ความหมายของระบบตรวจจับการบุกรุก

ระบบตรวจจับการบุกรุก คือ ระบบตรวจจับสัญญาณของความผิดปกติต่างๆ ที่เกิดขึ้นในระบบที่อยู่ในขอบเขตที่ระบบนี้มีหน้าที่ตรวจสอบ ในที่นี้จะหมายถึง โปรแกรมที่ใช้สำหรับตรวจจับความผิดปกติในระบบเครือข่ายคอมพิวเตอร์เท่านั้น โดยตัวโปรแกรมจะมีความสามารถในการตรวจจับสัญญาณของความผิดปกติที่เกิดขึ้นในระบบ ไม่ว่าจะเป็นภายในระบบคอมพิวเตอร์ ระบบปฏิบัติการ โปรแกรมที่รันอยู่ในเครื่อง การทำงานกับฐานข้อมูล หรือแม้แต่ข้อมูลที่วิ่งผ่าน ไปมาในเครือข่ายด้วย

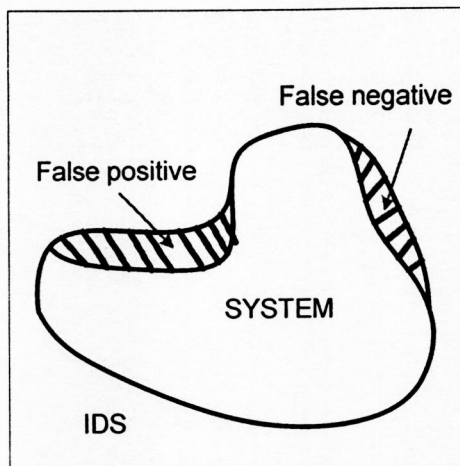


รูปที่ 2.1 ระบบการตรวจจับการบุกรุก

จากรูปที่ 2.1 ถ้าเรามองระบบเป็นชุดของการทำงานชุดหนึ่ง ระบบตรวจจับการบุกรุกที่แท้จริง (Ideal IDS) ต้องทราบขอบเขตของระบบว่าระบบทำงานอะไรบ้าง การทำงานใดปกติและการทำงานใดผิดปกติ โดยระบบตรวจจับการบุกรุกที่มีประสิทธิภาพ จะทราบขอบเขตของระบบ โดยไม่มีการเหลื่อมล้ำเข้าไปในระบบ หรือเหลื่อมล้ำออกนอกระบบ 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

การออกแบบระบบ IDS พยายามออกแบบเพื่อทำให้ False positive และ False negative มีน้อยที่สุดซึ่งก็มีงานวิจัยหลายๆ ชิ้นที่ทำเพื่อลดพื้นที่ตรงส่วนนี้ [1], [2], [3], [4], [5], [6]



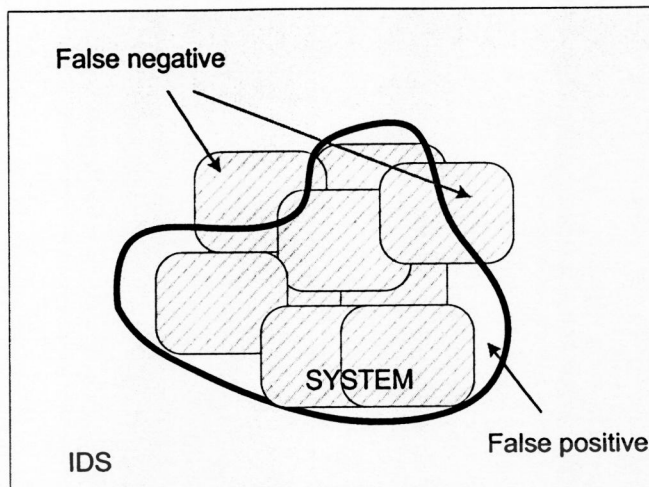
รูปที่ 2.2 ขอบเขตที่เหลื่อมกันของ IDS กับระบบ ทำให้เกิด False positive และ False negative

2.1.2 ประเภทของระบบตรวจจับการบุกรุก

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

2.1.2.1 Anomaly Detection

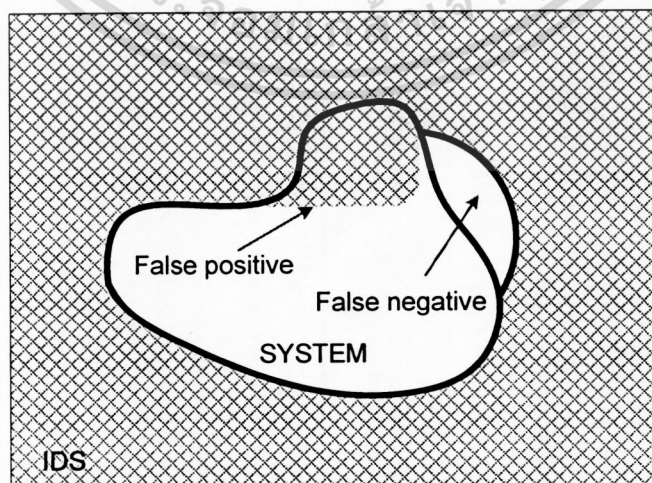
แนวความคิดของการทำ Anomaly Detection คือ การหาเซตของการทำงานที่เป็นปกติย่อยๆ ขึ้นมาแล้วนำมารวมกันเพื่อให้ระบบ IDS ทราบข้อมูลของเซตการทำงานที่เป็นปกติทั้งหมดในระบบ หลังจากนั้นเมื่อให้ระบบ IDS ทำงาน ถ้าเกิดกรณีที่ IDS ตรวจจับการทำงานที่ไม่ได้อยู่ในเซตของการทำงานที่เป็นปกติ ระบบ IDS จะแจ้งเตือนต่อผู้ดูแลระบบทันที 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

2.1.2.2 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



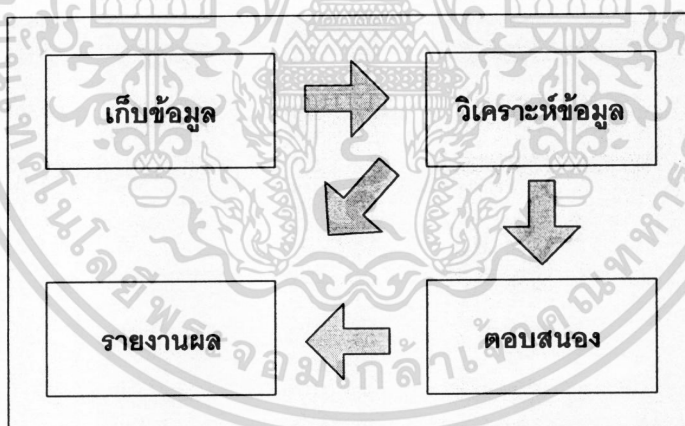
รูปที่ 2.4 Misuse Detection Model

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.1.3 การทำงานของระบบตรวจจับการบุกรุก

ระบบตรวจจับการบุกรุกแต่ละแบบมีหน้าที่การทำงานที่แตกต่างกันออกไป บางตัวจะตรวจจับความผิดปกติในระบบเครือข่าย บางตัวจะตรวจจับความผิดปกติในระบบฐานข้อมูล แต่โดยการทำงานทั้งหมดแล้วเราสามารถแบ่งการทำงานของระบบตรวจจับการบุกรุกได้เป็น 3 ขั้นตอน คือ การเก็บข้อมูลระบบ การวิเคราะห์ข้อมูลที่เก็บได้ และการรายงานผลการทำงานให้ผู้ดูแลระบบหรือผู้ที่เกี่ยวข้องทราบ

ในการทำงานหลักๆ ของระบบตรวจจับการบุกรุก อาจมีขั้นตอนเสริมอยู่ขั้นตอนหนึ่ง คือ การตอบสนองต่อการบุกรุกนั้นๆ การทำงานในขั้นตอนนี้จะใช้ในกรณีที่การบุกรุกเป็นรูปแบบการบุกรุกที่ระบบตรวจจับการบุกรุกสามารถแก้ไขด้วยตัวเองได้ ซึ่งในบางระบบอาจไม่มีการทำงานในส่วนนี้ ระบบที่ไม่มีการตอบสนองต่อการบุกรุก ส่วนใหญ่จะเป็นระบบที่ไม่ได้ทำงานแบบตอบสนองทันที (real time) คือ จะเก็บข้อมูลของระบบไว้ก่อน แล้วจึงวิเคราะห์ข้อมูลภายหลัง เมื่อทำการวิเคราะห์ข้อมูลแล้วพบว่ามีกรบุกรุกเข้าสู่ระบบก็จะทำการแจ้งเตือนในขั้นตอนการทำการรายงานผลการทำงาน ระบบตรวจจับการบุกรุกที่ไม่มีการตอบสนองต่อการบุกรุกก็มักใช้ในงานที่ไม่มีความสำคัญมากนัก แต่ต้องการความถูกต้องสูง โดยลำดับการทำงานของระบบตรวจจับการบุกรุกสามารถมองเป็นขั้นตอนต่างๆ ได้ดังรูปที่ 2.5



รูปที่ 2.5 การทำงานของระบบตรวจจับการบุกรุก

2.1.3.1 การเก็บข้อมูลในระบบ

จากที่ได้กล่าวมาแล้วว่าระบบตรวจจับการบุกรุกจะมีการทำงานที่แตกต่างกันไป หน้าที่ในการเก็บข้อมูลของระบบที่ต้องการตรวจสอบก็แตกต่างกันไปตามหน้าที่ของระบบตรวจจับผู้บุกรุกด้วย โดยเราสามารถแบ่งการเก็บข้อมูลของระบบที่ต้องการตรวจสอบออกเป็นกลุ่มต่างๆ ได้ 4 กลุ่มด้วยกันคือ มีการเก็บข้อมูลในชั้นแอปพลิเคชัน (Application-based Approach) เพื่อ

นำมาตรวจสอบการทำงานของแอปพลิเคชันต่างๆ ว่าผิดปกติหรือไม่ การเก็บข้อมูลของการดำเนินงานของเครื่อง (Host-based Approach) เพื่อนำมาตรวจสอบการทำงานของระบบของเครื่องที่ใช้

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

2.1.3.1.1 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

2.1.3.1.2 การเก็บข้อมูลของการทำงานของเครื่อง

สำหรับการเก็บข้อมูลของการทำงานของเครื่อง จะเน้นไปในการเก็บข้อมูลของระบบปฏิบัติการเป็นหลัก ข้อมูลที่เก็บได้จะอยู่ในรูปของการแจ้งเตือนในระบบเช่นการตั้งค่าบางอย่างไม่สมบูรณ์ การทำงานของโพรเซสบางโพรเซสมีปัญหาหรือปัญหาของฮาร์ดแวร์เป็นต้น หรืออาจอยู่ในรูปข้อมูลของการทำงานของระบบปฏิบัติการนั้นๆ เช่น ข้อมูลการใช้งานของยูสเซอร์แต่ละคน ใคร ทำอะไร เมื่อเวลาเท่าไร ซึ่งเมื่อนำข้อมูลเหล่านี้ไปวิเคราะห์แล้ว จะได้ผลการวิเคราะห์ในลักษณะมีการใช้งานอย่างไม่ถูกต้องหรือไม่ ถ้ามี ใครเป็นผู้ใช้งานนั้นๆ เมื่อเวลาเท่าไรจากที่ไหน ข้อดีอีกข้อหนึ่งก็คือการเก็บข้อมูลในลักษณะนี้สามารถเก็บข้อมูลที่ถูกต้องเข้ารหัสได้ด้วยส่วนข้อเสียของการเก็บข้อมูลการทำงานของเครื่องก็คือ ข้อมูลที่ได้มักจะมีขนาดใหญ่ ระบบที่ทำการเก็บข้อมูลลักษณะนี้จะมี Overhead สูงขึ้น นอกจากนี้โปรแกรมที่ทำการเก็บข้อมูลและวิเคราะห์ข้อมูลยังขึ้นอยู่กับ platform และมีราคาสูงมากด้วย

2.1.3.1.3 การเก็บข้อมูลการเปลี่ยนแปลงข้อมูลในระบบ

การเก็บข้อมูลการเปลี่ยนแปลงข้อมูลในระบบจะใช้หลักการของ integrity analysis ในการตรวจสอบการเปลี่ยนแปลงข้อมูลต่างๆ ในระบบ ในระบบตรวจจับการบุกรุกบางระบบจะใช้ checksum เป็นตัวบ่งบอกการเปลี่ยนแปลงในระบบ การวิเคราะห์ลักษณะนี้จะเริ่มจากการสร้างฐานข้อมูล signature ของไฟล์ต่างๆ ในระบบปกติไว้ ในระบบปกติไว้ เมื่อระบบมีการทำงานก็จะทำการตรวจสอบค่า signature นี้ไปเรื่อยๆ อาจเป็นวันละครั้ง สองครั้ง หรือบ่อยกว่านั้นแล้วแต่

เอกสารนี้เป็นเอกสารที่สงวนไว้เพื่อใช้ในการเรียนการสอนเท่านั้น ไม่สามารถนำไปเผยแพร่หรือใช้เพื่อการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

integrity analysis ลักษณะนี้ช่วยให้การตรวจจับการบุกรุกที่มีการเปลี่ยนแปลงระบบ เช่น ทำการเจาะระบบแล้วทำการวางโทรจัน หรือ Back Door ไว้ได้ สำหรับการแก้ไขเมื่อทราบว่าการบุกรุกมาเปลี่ยนแปลงไฟล์ข้อมูลในระบบก็ทำการแก้ไขเฉพาะไฟล์ที่ถูกแก้ไขเท่านั้น ไม่จำเป็นต้องทำการติดตั้งระบบใหม่ แต่ระบบนี้ก็ยังมีข้อเสียในกรณีที่มีไฟล์ในระบบเป็นจำนวนมาก การเก็บข้อมูล signature ของไฟล์ต่างๆ และการวิเคราะห์ข้อมูลก็จะใช้เวลานาน ระบบนี้จึงไม่เหมาะในการทำงานแบบทันที (real time) เพราะจะทำให้เกิด overhead ในระบบสูงมาก

2.1.3.1.4 การเก็บข้อมูลเครือข่าย

การเก็บข้อมูลเครือข่ายนั้นนับวันจะมีความสำคัญขึ้นเรื่อยๆ เพราะการบุกรุกทางเครือข่ายมีมากขึ้นเรื่อยๆ การเก็บข้อมูลแบบนี้จะใช้การดักจับข้อมูลที่ผ่านไปมาในเครือข่าย โดยการทำให้เน็ตเวิร์คการ์ดอยู่ใน promiscuous mode เมื่อเน็ตเวิร์คการ์ดอยู่ในโหมดดังกล่าว จะสามารถรับข้อมูลทุกอย่างที่อยู่ในเครือข่ายได้ การเก็บข้อมูลเครือข่ายในลักษณะนี้สามารถตรวจจับการโจมตีทางเครือข่ายได้ เช่น การทำ SYN flood การทำ port scan หรือการส่งแพ็กเก็ตปริมาณมากมารบกวนในระบบ แต่เนื่องจากการเก็บข้อมูลเครือข่ายนี้ใช้ลักษณะการนำสนิฟเฟอ์เป็นหลัก จึงไม่สามารถทำงานในเครือข่ายที่เป็นเครือข่ายสวิตซ์ซึ่ง ไม่สามารถทำงานในระบบเครือข่ายที่เข้ารหัสข้อมูล หรือไม่สามารเก็บข้อมูลในเครือข่ายที่มีข้อมูลหนาแน่น ได้เพราะการทำงานในการเก็บข้อมูลอาจไม่เร็วพอที่จะเก็บข้อมูลทั้งหมดที่ผ่านไปมาในระบบได้ ข้อเสียอีกข้อหนึ่งของการเก็บข้อมูลนี้ก็คือข้อมูลที่เก็บมีขนาดใหญ่มาก โดยเฉพาะอย่างยิ่งในระบบเครือข่ายที่มีการรับส่งแพ็กเก็ตปริมาณมากอยู่ตลอดเวลา

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

2.1.3.2 การวิเคราะห์ข้อมูลระบบ

เมื่อได้ข้อมูลของระบบที่จำเป็นแล้ว ในขั้นตอนต่อมาเราก็จะนำเอาข้อมูลที่ได้มาวิเคราะห์ว่าระบบของเรามีความผิดปกติเกิดขึ้นหรือไม่ การวิเคราะห์ข้อมูลสามารถแบ่งการทำงานตามรูปแบบการวิเคราะห์ข้อมูลได้ 2 รูปแบบ คือ ทำการวิเคราะห์ในขณะที่เก็บข้อมูลแบบทันทีหรือจะเก็บข้อมูลทั้งหมดไว้ก่อน แล้วจึงวิเคราะห์ข้อมูลนั้นๆ ภายหลัง (Batch) ในการวิเคราะห์ข้อมูลทั้งสองรูปแบบก็มีข้อเสียแตกต่างกันไป

2.1.3.2.1 การวิเคราะห์ในขณะที่เก็บข้อมูลทันที (Real Time)

ในการวิเคราะห์ข้อมูลที่ได้ในขณะที่เก็บข้อมูล หรือแบบทันทีนั้น ระบบจะจัดเก็บข้อมูล

วิเคราะห์ข้อมูล และรายงานผลการวิเคราะห์ ในช่วงเวลาเดียวกัน เมื่อเกิดข้อผิดพลาดขึ้นสามารถ
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตอบสนองได้ทันทีทั้งที่ระบบที่ทำงานแบบทันทีมีการแจ้งเตือนหลายๆ แบบ เช่น E-mail หรือ Instant Messaging ให้กับผู้ดูแลระบบได้ในเวลาที่มีการบุกรุกได้ ในการตรวจสอบระบบแบบทันทีทำให้ระบบสามารถตรวจสอบข้อผิดพลาดได้อย่างรวดเร็ว แต่ก็ขึ้นอยู่กับความเร็วในการวิเคราะห์ข้อมูลด้วย ถ้าข้อมูลมีความซับซ้อนมากๆ ก็จะใช้เวลามากตามเช่นเดียวกันเมื่อระบบทำการตรวจสอบการบุกรุกได้ในขณะที่เพิ่งเกิดการบุกรุกขึ้น ผู้ดูแลระบบ หรือ ระบบตรวจจับการบุกรุกเองสามารถแก้ไขปัญหาที่เกิดขึ้นได้ทันที แต่ทั้งนี้ก็ขึ้นอยู่กับความเร็วในการวิเคราะห์ข้อมูล และชนิดของปัญหาว่ามีความยุ่งยากในการแก้ปัญหาเล็กน้อยเพียงไรด้วย

ระบบวิเคราะห์ข้อมูลแบบทันที ทำงานได้อย่างรวดเร็ว แต่การทำงานที่รวดเร็วดังกล่าวก็ ต้องแลกกับการใช้หน่วยความจำปริมาณมาก และการประมวลผลที่รวดเร็วมักด้วย อีกทั้งการตอบสนองต่อการบุกรุกโดยอัตโนมัติ อาจทำให้เกิดความเสียหายกับระบบมากกว่าเดิม เพราะในบางครั้ง การทำงานที่เร็วเกินไปของระบบนี้ ทำให้เกิดความผิดพลาดในการวิเคราะห์ข้อมูลจนประมวลผลการทำงานปกติ กลายเป็นการทำงานที่ผิดปกติ แล้วทำการแก้ไขตามข้อมูลที่มีอยู่ ก็ยังทำให้ระบบมีความเสียหายมากกว่าเดิม ระบบตรวจจับการบุกรุกที่ทำงานแบบทันทีจึงเหมาะกับระบบที่มีข้อมูลที่ต้องพิจารณาบ่อย ต้องการการรายงานอย่างรวดเร็วเมื่อผิดปกติ และข้อมูลที่ต้องนำมาวิเคราะห์ไม่ซับซ้อนมากนัก

2.1.3.2.2 การวิเคราะห์ข้อมูลภายหลังจากที่เก็บข้อมูล (Batch)

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

ไม่ว่าจะเป็นการวิเคราะห์ระบบแบบทันทีหรือ Batch ก็จะมีวิธีการวิเคราะห์ระบบที่เหมือนกัน คือ ทำการหารูปแบบของการโจมตีในข้อมูลที่ได้รับมา (Signature Analysis) วิธีการวิเคราะห์ทางสถิติ (Statistical Analysis) และวิธีการตรวจสอบการเปลี่ยนแปลงของระบบ (Integrity Analysis)

1) การหารูปแบบของการโจมตี (Signature Analysis)

วิธีการวิเคราะห์ระบบแบบ Signature Analysis เป็นการวิเคราะห์ข้อมูลโดยการหาสัญญาณของการโจมตี (Attack Signature) การทำงานจะทำโดยการเปรียบเทียบรูปแบบของข้อมูลกับรูปแบบของการโจมตีในฐานข้อมูล ว่ามีความคล้ายกันหรือไม่ ถ้ามีความคล้ายคลึงกันก็แสดงว่ามีการโจมตีเกิดขึ้นแล้ว ในการเปรียบเทียบอาจเป็นแบบอย่างง่าย คือ การเปรียบเทียบข้อมูลว่ามี

ความเข้ากันได้กับข้อมูลของการโจมตีเพียงใด 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 สำหรับโพรไฟล์นั้นอาจแยกเป็นเป็นโพรไฟล์สำหรับแอปพลิเคชันต่างๆ ในระบบ เช่น ยูสเซอร์ ไฟล์ ไคลเอนท์ และอุปกรณ์ต่างๆ โดยรวมละเอียดที่เก็บอยู่ในโพรไฟล์จะมีข้อมูลของจำนวนครั้งที่เข้าสู่ระบบ จำนวนครั้งที่เข้าสู่ระบบผิดพลาด เวลา และข้อมูลอื่นๆ ที่จำเป็น ค่าแต่ละค่าที่เก็บจะเป็นค่าของการใช้งานที่เป็นปกติ การตรวจจับว่าเกิดความผิดปกติขึ้นแล้ว จะดูจากค่าที่ไม่เข้ากับโพรไฟล์ เป็นต้น ยกตัวอย่างเช่น ในการทำงานปกติ ผู้ใช้งานฐานข้อมูลจะมีการเข้าใช้ข้อมูลในฐานข้อมูลตั้งแต่เวลา 8 โมงเช้า ถึง 6 โมงเย็นเท่านั้น แต่ข้อมูลที่ตรวจจับได้มีการเข้าใช้ฐานข้อมูลตอนตีสอง ซึ่งก็สามารถบอกได้ว่าการบุกรุกเกิดขึ้นแล้ว

เนื่องจากวิธีการวิเคราะห์ข้อมูลทางสถิติ เป็นการตรวจจับโดยใช้หลักการของการอนุญาตให้ใช้งานในการทำงานต่างๆ ไปและไม่อนุญาตให้ใช้งานนอกเหนือจากที่เคยใช้โดยทั่วไปเท่านั้น การตรวจจับในลักษณะนี้จึงสามารถตรวจจับการบุกรุกที่ไม่เคยเจอมาก่อนได้ และสามารถตรวจจับการบุกรุกในรูปแบบที่ซับซ้อนได้ด้วย เพราะเราถือว่าการทำงานที่สลับซับซ้อนมักจะไม่เหมือนกับการทำงานโดยปกติ แต่วิธีการวิเคราะห์ข้อมูลแบบนี้ก็มีข้อเสียเนื่องจากการเก็บข้อมูลการทำงานที่เป็นปกติไว้เพื่อเปรียบเทียบกับการทำงานที่ผิดปกติ เมื่อทำการบุกรุกในลักษณะเดิมๆ เป็นเวลานาน ก็จะทำให้โพรไฟล์มีการเปลี่ยนแปลง ระบบตรวจจับก็จะเห็นว่าการโจมตีในลักษณะนั้นเป็นการทำงานที่เป็นปกติแทน และไม่สามารถตรวจจับการทำงานที่ผิดปกติในลักษณะนั้นได้อีกต่อไป และไม่เหมาะกับองค์กรที่มีการเปลี่ยนแปลงการทำงานบ่อยๆ เพราะ

เอกส ทำให้โพรไฟล์มีขนาดใหญ่ ทำให้ระบบตรวจจับรวน และมีความผิดพลาดในการวิเคราะห์สูง ด้านการคำนวณว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3) วิธีการตรวจสอบการเปลี่ยนแปลงของระบบ (Integrity Analysis)

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

2.1.3.3 การตอบสนอง

เมื่อมีการตรวจพบว่ามีกรบุกรุกเกิดขึ้นในระบบ สำหรับระบบตรวจจับที่ทำงานแบบ real time จะมีการตอบสนองต่อการบุกรุกเพื่อไม่ให้เกิดความเสียหาย หรือบรรเทาความเสียหายที่เกิดขึ้นสำหรับระบบที่ทำงานเป็นแบบ batch การตอบสนองอาจทำได้ไม่มากนัก เพราะการบุกรุกนั้นเกิดไปแล้ว ความเสียหายก็เกิดขึ้นแล้ว การตอบสนองอาจอยู่ในรูปแบบการบรรเทาไม่ให้ความเสียหายมีเพิ่มมากขึ้นเท่านั้น การตอบสนองต่อการบุกรุกนั้นแบ่งออกได้เป็นสามแบบด้วยกัน คือ การเปลี่ยนแปลงสภาพของระบบ การแก้ไขความผิดพลาดให้ถูก และการแจ้งเตือนผู้ดูแลระบบเมื่อถูกบุกรุก

2.1.3.3.1 การเปลี่ยนแปลงสภาพของระบบ

สำหรับการตอบสนองต่อการบุกรุกโดยการเปลี่ยนแปลงสภาพของระบบที่ถูกโจมตีก็เพื่อแก้ปัญหาหรือลดความเสียหายที่จะเกิดขึ้น เช่น ตัดการเชื่อมต่อระหว่างระบบกับการบุกรุกออกจากกัน การตั้งค่าอุปกรณ์เครือข่าย หรือไฟร์วอลล์ไม่ให้มีการติดต่อกับระบบของการบุกรุกอีกต่อไป และการหาข้อมูลเกี่ยวกับการโจมตีโดยอัตโนมัติเพื่อตรวจหาการบุกรุกต่อไป

2.1.3.3.2 การแก้ไขความผิดพลาดให้ถูก

การแก้ไขระบบ เป็นการตอบสนองต่อปัญหาที่เกิดขึ้นแล้วในระบบ โดยปกติแล้วการบุกรุกมักเปลี่ยนแปลงค่าต่างๆ ในระบบ โดยเฉพาะเข้ามาทำการเปลี่ยนแปลงข้อมูลในระบบตรวจจับการบุกรุกเพื่อไม่ให้อาจสามารถตรวจจับการบุกรุกได้ การแก้ไขระบบก็เพื่อให้ระบบดังกล่าวสามารถทำงานได้อย่างเป็นปกติ

2.1.3.3.3 การแจ้งเตือนผู้ดูแลระบบ

สุดท้ายเป็นการแจ้งเตือนผู้ดูแลระบบ โดยปกติมักแจ้งเตือนผู้ดูแลทันทีเมื่อทำการวิเคราะห์ได้ว่ามีความผิดปกติเกิดขึ้น เพื่อให้ผู้ดูแลระบบรับรู้และสามารถแก้ไขระบบได้ทันเวลาที่ สำหรับการแจ้งเตือนนี้ ผู้ดูแลระบบสามารถเลือกได้ว่าจะแจ้งเตือนใครบ้าง และทำการแจ้งเตือนในรูปแบบใด เช่น E-mail, Pager หรือ Instant Messaging เป็นต้น

2.1.3.4 การรายงานผลการทำงาน

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

2.1.4 ความสำคัญของระบบตรวจจับการบุกรุก

เมื่อได้ทราบถึงการทำงานคร่าวๆ ของระบบตรวจจับการบุกรุกแล้ว อาจคิดว่าระบบตรวจจับการบุกรุกไม่มีความสำคัญเพราะในเมื่อมีการใช้งานไฟร์วอลล์อยู่แล้ว แต่ความเป็นจริงแล้ว ถึงแม้ว่าระบบจะมีไฟร์วอลล์อยู่แล้วก็ยังจำเป็นต้องใช้ระบบตรวจจับการบุกรุกด้วยเพราะในบางอย่างไฟร์วอลล์ก็ไม่สามารถช่วยได้

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

2.2 งานวิจัยที่เกี่ยวข้อง

ปัจจุบันการวิจัยที่เกี่ยวกับระบบตรวจจับการบุกรุกได้เป็นหัวข้อวิจัยที่แพร่หลาย เนื่องจากมีการบุกรุกระบบเครือข่ายคอมพิวเตอร์เพิ่มมากขึ้นตามจำนวนคอมพิวเตอร์ที่ต่อเข้าสู่เครือข่าย โดยนักวิจัยได้พัฒนาการตรวจจับการบุกรุกอย่างต่อเนื่องโดยใช้วิธีต่างๆ ดังนี้

Wenke Lee และ Salvatore J. Stolfo [10] ได้เสนอวิธีตรวจจับการบุกรุกโดยใช้เทคนิคคาด้า ไม่นิ่งคำนวณหารูปแบบของพฤติกรรมจากข้อมูลที่ตรวจสอบจากระบบและดึงคุณลักษณะที่ใช้จำแนกพฤติกรรมการบุกรุกออกจากรูปแบบพฤติกรรม จากนั้นใช้เทคนิคการเรียนรู้ของเครื่องจักรสร้างกฎตามนิยามของคุณลักษณะ ผลการทดลองกับข้อมูลที่ตรวจสอบที่ได้จากโปรแกรม the 1998 DARPA Intrusion Detection Evaluation Program ได้ค่า Detection rate เท่ากับ 80.2%

Mahdireza Mohajerani Ali Moeini และ Mojtaba Kianie [11] ได้เสนอระบบตรวจจับ NFIDS (The Neuro-Fuzzy Intrusion Detection System) ซึ่งเป็นระบบตรวจจับการบุกรุกแบบ anomaly ที่มีสถาปัตยกรรมแบบ hierachical ที่ใช้พีชชี่ลอจิกและนิเวรอลเน็ตเวิร์คตรวจจับพฤติกรรมการบุกรุกบนเครือข่าย โดยใช้นิเวรอลเน็ตเวิร์คแบบ MLP (Multi-Layer Perceptron) เรียนรู้กฎพีชชี่ที่สร้างจากการบุกรุกแบบพอร์ตสแกนและทำการทดสอบนิเวรอลเน็ตเวิร์คกับข้อมูลที่อยู่บนเครือข่ายจริงๆ ผลการทดลองได้ค่า Detection rate เท่ากับ 90.6% False negative rate เท่ากับ 4% และ False positive rate เท่ากับ 9.4%

Jonatan Gomez และ Dipankar Dasgupta [12] ได้นำเจนเนติกอัลกอริทึมมาสร้างกฎพีชชี่ที่ใช้จำแนกการบุกรุกที่สามารถจำแนกการบุกรุกในแบบ Anomaly และการบุกรุกที่กำหนด โดยได้ทำการแบ่งกฎพีชชี่ที่ใช้จำแนกการบุกรุกที่สร้างขึ้นออกเป็นสองประเภท คือกฎพีชชี่ที่ใช้จำแนกพฤติกรรมปกติและกฎพีชชี่ที่ใช้จำแนกพฤติกรรมบุกรุก โดยใช้ชุดข้อมูลจากฐานข้อมูลของ KDD Cup 1999 ที่มีรูปแบบพฤติกรรมของระบบเป็นแบบปกติและบุกรุก จากการทดลองกฎพีชชี่ที่ใช้จำแนกการบุกรุกที่สร้างขึ้นด้วยฐานข้อมูลของ KDD Cup 1999 พบว่าจะได้ค่า Detection rate เท่ากับ 95.47% และ False positive rate เท่ากับ 10.63%

Chaivat Jirapummin Naruemon Wattanapongsakorn และ Prasert Kanthamanon [5] ได้นำเสนอระบบตรวจจับการบุกรุกที่ใช้ Self-Organizing Maps (SOM) และ Resilient Propagation Neural Network (RPROP) ในการจำแนกพฤติกรรมปกติและพฤติกรรมบุกรุก โดยใช้ SOM ในการจัดกลุ่มของพฤติกรรมปกติและพฤติกรรมบุกรุก โดยแสดงการจัดกลุ่มอยู่ในรูปแบบแผนภาพสองมิติ และใช้ RPROP จำแนกพฤติกรรมปกติและพฤติกรรมบุกรุก โดยใช้ชุดข้อมูลจากฐานข้อมูลของ KDD Cup 1999 สำหรับทำการสอนและทดสอบระบบ จากการทดลองระบบด้วยฐานข้อมูลของ KDD Cup 1999 พบว่าได้ค่า Detection rate มากกว่า 90% และ False alarm rate

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

น้อยกว่า 5% ในการจำแนกการบุกรุกแบบ SYN flooding หนึ่งประเภท คือ neptune และการบุกรุกแบบ port scanning สองประเภท คือ port sweep และ satan

Zhi-Song Pan Song-Can Chen Gen-Bao Hu และ Dao-Qiang Zhang [13] ได้นำเสนอโมเดลตรวจจับการบุกรุกแบบ misuse ที่ใช้ Hybrid Neural Network และ C4.5 โดยนำความสามารถในการจำแนกพฤติกรรมกรการบุกรุกที่ไม่เหมือนกันของ Neural Network และ C4.5 Algorithm มาใช้ในการออกแบบโมเดล โดยใช้ชุดข้อมูลจากฐานข้อมูลของ KDD Cup 1999 สำหรับทำการสอนและทดสอบระบบ อันดับแรกจะทำการสอนและทดสอบ hybrid model ด้วยข้อมูลของพฤติกรรมปกติและพฤติกรรมกรการบุกรุกที่รู้จัก จากนั้นจะทำการปรับแก้ hybrid model ให้มีประสิทธิภาพในการจำแนกพฤติกรรมกรการบุกรุกดีขึ้น เนื่องจาก BP network และ C4.5 มีความสามารถในการจำแนกพฤติกรรมกรการบุกรุกไม่เหมือนกัน จากการทดลองโมเดลด้วยฐานข้อมูลของ KDD Cup 1999 พบว่าจะได้ค่า Detection rate เฉลี่ยเท่ากับ 85% และ False alarm rate สำหรับจำแนกพฤติกรรมกรการบุกรุก 5 ประเภท คือ neptune , portsweep , satan , buffer_overflow , และ guess_passwd น้อยกว่า 19.7% หลังจากแก้ไขโมเดลให้ดีขึ้น จะได้ค่า Detection rate เฉลี่ยเท่ากับ 93.28% และ False positive rate เท่ากับ 0.2%

2.3 สรุป

ระบบตรวจจับการบุกรุกเป็นผู้ช่วยที่ดีสำหรับผู้ดูแลระบบ หน้าที่ของการตรวจจับการบุกรุกนั้นจะรวบรวมข้อมูล วิเคราะห์ข้อมูลและทำการแจ้งผลการวิเคราะห์ข้อมูลต่างๆ ในระบบให้กับผู้ดูแลระบบ ซึ่งช่วยให้การดูแลระบบทำได้อย่างมีประสิทธิภาพ แต่ระบบตรวจจับการบุกรุกก็ยังไม่ใช่ว่าสิ่งที่จะมาแก้ปัญหาคความปลอดภัยในระบบโดยสิ้นเชิงได้ เนื่องจากการที่จะทำให้ระบบปลอดภัยต้องอาศัยความร่วมมือจากหลายๆ ฝ่าย ไม่เพียงแต่เฉพาะการดูแลของผู้ดูแลระบบคนเดียว สิ่งที่สำคัญที่สุดที่ผู้เขียนเห็นว่าจะสร้างความปลอดภัยให้กับระบบในระยะยาวก็คือ การปลูกจิตสำนึกด้านความปลอดภัยในการใช้งานให้กับผู้ใช้ระบบแต่ละคนด้วย

งานวิจัยที่กล่าวถึงในหัวข้อก่อนหน้านี ยังมีค่า detection rate ต่ำ และ false alarm rate ที่สูง นอกจากนี้กฎที่ใช้ในการจำแนกการบุกรุกที่สร้างยังไม่สามารถอธิบายได้ เพื่อแก้ปัญหาดังกล่าว งานวิจัยนี้จึงได้เสนอระบบตรวจจับการบุกรุกผ่านเครือข่ายแบบ Misuse โดยใช้ราฟฟิชชี ซึ่งจะกล่าวรายละเอียดในบทที่ 5 ต่อไป

บทที่ 3

ฟัซซีเซต

ทฤษฎีฟัซซีเซตได้ถูกพัฒนาขึ้นโดย Zadeh ในปี ค.ศ. 1965 ซึ่งมีความพยายามที่จะเลียนแบบหลักการคิดของมนุษย์ที่สามารถพิจารณาแก้ไขปัญหาต่างๆ ได้โดยใช้ความรู้ เหตุผล และปัจจัยแวดล้อมต่างๆ มาประกอบกัน เพื่อหาทางออกของปัญหาที่เหมาะสมที่สุดและถูกนำมาใช้ในการอธิบายระบบที่มีความคลุมเครือ ฟัซซีเซตมีพื้นฐานมาจากทฤษฎีเซต โดยค่าความเป็นสมาชิกของข้อมูลจะอยู่ในช่วง 0 ถึง 1 ซึ่งมีความแตกต่างจากค่าความเป็นสมาชิกของเซตธรรมดาที่มีค่า 0 และ 1 เท่านั้น ทฤษฎีของฟัซซีเซตได้ถูกพัฒนาไปใช้งานหลายๆ ด้าน เช่น งานทางด้านระบบควบคุมในทางวิศวกรรมศาสตร์ ทางด้านประมวลผลภาพหรือการจัดกลุ่มของข้อมูล เป็นต้น

3.1 ทฤษฎีฟัซซีเซต

ในระบบที่เป็นคริสป์ (crisp set) การระบุถึงการเป็นสมาชิกของเซตจะมีเพียง “เป็นสมาชิก” และ “ไม่เป็นสมาชิก” โดยแทนด้วยค่าความเป็นสมาชิก 0 หรือ 1 แต่ในฟัซซีเซต การระบุถึงการเป็นสมาชิกจะอ้างจากค่าความเป็นสมาชิกที่มีค่าอยู่ในช่วง 0 ถึง 1 ถ้ากำหนดให้ U เป็นเซตเอกภพสัมพัทธ์ และฟัซซีเซต $A \subseteq U$ ดังนั้นฟัซซีเซต A สามารถเขียนให้อยู่ในรูปความสัมพันธ์ของ $x \in U$ กับ $\mu_A(x)$ ดังสมการที่ 3.1

$$A = \{(x, \mu_A(x)) | x \in U\} \tag{3.1}$$

โดยที่ $\mu_A(x)$ คือฟังก์ชันความเป็นสมาชิกของ x ในฟัซซีเซต A

ตัวแปรฟัซซี (Fuzzy variable) หรือบางครั้งอาจจะเรียกว่า Linguistic variable คือฟัซซีเซตของระบบที่เราสนใจ ตัวอย่างเช่น ถ้าเรากำหนดให้ “ความสูง” เป็นตัวแปรฟัซซี และกำหนดเซตค่าของตัวแปรเป็น {เตี้ย, ปานกลาง, สูง} เราอาจจะเรียกเซตของค่าตัวแปรเหล่านี้ว่าค่าตัวแปรฟัซซี หรือเทอมเซต (Term set) นอกจากนี้ตัวแปรฟัซซีแต่ละตัวอาจจะมีส่วนขยาย (Hedges) เพื่อปรับค่าตัวแปรฟัซซีให้มีความยืดหยุ่นมากยิ่งขึ้น เช่น จากค่าของตัวแปรความสูง “เตี้ย-เตี้ยมาก” หรือ “สูง-ค่อนข้างสูง” เป็นต้น คำว่า “มาก” หรือ “ค่อนข้าง” ในที่นี้เป็นส่วนขยายของตัวแปรฟัซซี

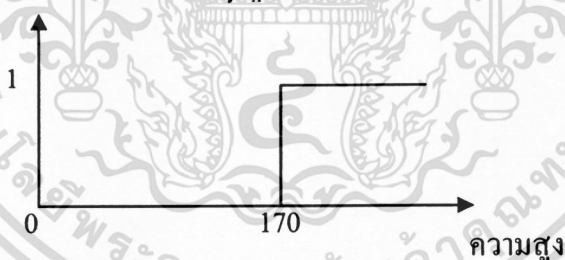
3.2 การแทนข้อมูลในฟัซซีเซต

การแทนข้อมูลในฟัซซีเซต เพื่อให้เข้าใจมากยิ่งขึ้นจะเปรียบเทียบกับระบบเซตธรรมดา ตัวอย่างเช่น “ความสูง” (ในระบบฟัซซีเซต ก็คือตัวแปรฟัซซีเซต) ที่เราสนใจอยู่ในช่วง 140-180 เซนติเมตร (เป็นโดเมนของเอกภพสัมพัทธ์) แต่ถ้าเราสนใจเฉพาะคนที่สูง และกำหนดให้เซต A เป็นเซตของคนที่สูง (“สูง” ในฟัซซีเซตจะเป็นเทอมเซต หรือค่าตัวแปรฟัซซี) ดังนั้นเซต A ในกรณีของระบบเซตธรรมดาจะต้องมีการกำหนดค่าเริ่มเปลี่ยน หรือค่าเทรชโฮลด์ (Threshold) เพื่อเป็นตัวชี้ว่าข้อมูลจะเป็นสมาชิกของเซตใด ถ้าวัดสมมติกำหนดค่าเทรชโฮลด์มีค่าเท่ากับ 170 เซนติเมตร นั่นคือ ถ้าวัดใครมีความสูงมากกว่าหรือเท่ากับ 170 เซนติเมตร จะถือว่าคนนั้นสูง (เป็นสมาชิกของเซต A) ซึ่งสามารถเขียนเป็นฟังก์ชันความเป็นสมาชิกได้ดังสมการที่ 3.2

$$\mu_A(x) = \begin{cases} 1 & x \geq 170 \\ 0 & \text{otherwise} \end{cases} \quad (3.2)$$

โดยที่ x คือความสูงของคน ดังนั้นจากสมการที่ 3.2 เราสามารถนำมาเขียนเป็นกราฟแสดงความ เป็นสมาชิกของเซตได้รูปที่ 3.1

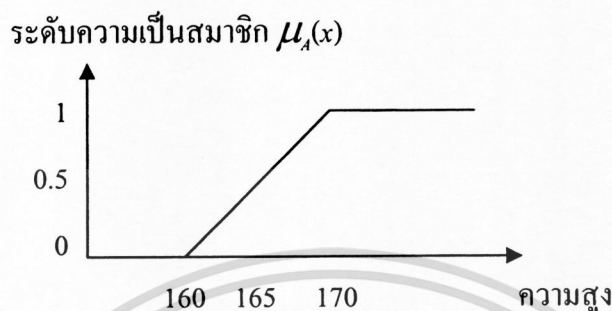
ระดับความเป็นสมาชิก $\mu_A(x)$



รูปที่ 3.1 กราฟแสดงค่าความเป็นสมาชิกของ x ที่มีต่อเซต A กรณีที่เป็นคริสป์เซต (crisp set)

จากรูปที่ 3.1 จะเห็นได้ว่าคนที่มีความสูงตั้งแต่ 170 เซนติเมตร ขึ้นไปเท่านั้นจึงถือว่าเป็นคนสูง (มีค่าความเป็นสมาชิกเท่ากับ 1) นอกนั้นถือว่าเป็นคนเตี้ยทั้งหมด แต่ในความเป็นจริงจะไม่ถูกต้องมากนัก เพราะคนที่มีความสูง 169.9 เซนติเมตร ก็มีความสูงใกล้เคียงกับคนที่สูง 170 เซนติเมตรมาก แต่ถือว่าเป็นคนเตี้ย จากปัญหานี้ สามารถอธิบายโดยฟัซซีเซตที่มีฟังก์ชันความเป็นสมาชิกที่มีความสอดคล้องกับความเป็นจริงมากกว่าการอธิบายโดยใช้คริสป์เซต ดังแสดงในสมการที่ 3.3

$$\mu_A(x) = \begin{cases} \frac{1}{10}(x-160) & 160 \leq x \leq 170 \\ 1 & x > 170 \\ 0 & \text{otherwise} \end{cases} \quad (3.3)$$

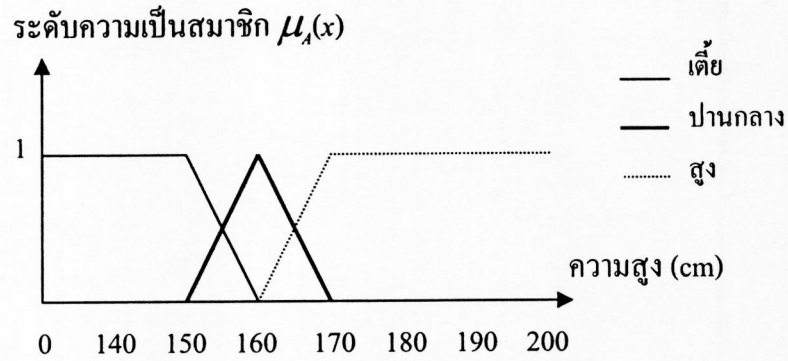


รูปที่ 3.2 กราฟแสดงค่าความเป็นสมาชิกของเซต A “ความสูง” กรณีที่เป็นฟuzzyเซต

จากรูปที่ 3.1 และรูปที่ 3.2 จะเห็นได้ว่ากราฟแสดงความเป็นสมาชิกของฟuzzyเซตจะมีความต่อเนื่องในช่วง $[0,1]$ เช่น คนที่มีความสูง 165 เซนติเมตร ก็จะถือว่าเป็นคนสูงเช่นกัน แต่เขาจะมีระดับความเป็นสมาชิกของเซตคนสูงเท่ากับ 0.5 แต่ถ้าเปรียบเทียบกับคริสป์เซตจะพบว่าคนที่มีความสูง 165 เซนติเมตรจะมีค่าความเป็นสมาชิกของเซตคนสูง (เซต A) เท่ากับค่าศูนย์ต่างๆ ที่ความสูงน้อยกว่า 170 เซนติเมตรเพียงเล็กน้อยเท่านั้น

3.3 ฟังก์ชันความเป็นสมาชิก

พิจารณาคูสมบัติของฟuzzyเซต จากกราฟรูปที่ 3.2 แกนนอนของกราฟจะแทนโดเมนของฟuzzyเซต ส่วนแกนตั้งจะแทนค่าระดับความเป็นสมาชิก แต่ในกราฟรูปที่ 3.2 มีเทอมเซตเพียงเทอมเดียว คือเทอมของคนสูงเท่านั้นเพื่อให้ครอบคลุมกลุ่มประชากรที่เราสนใจจึงเพิ่มเทอมเซตอีกสองเทอมคือ เทอมเซตของคน “เตี้ย” และเทอมเซตของคนที่มีความสูง “ปานกลาง” ดังรูปที่ 3.3 ซึ่งในกรณีที่โดเมนของตัวแปรฟuzzyเซตที่ครอบคลุมประชากรทั้งหมดที่ทำการศึกษา เราจะเรียกโดเมนของตัวแปรนี้ว่า เซตเอกภพสัมพัทธ์ และในแต่ละเทอมเซตก็จะมีโดเมนของตัวเองอย่างเช่น เทอมเซตของคนเตี้ยจะมีโดเมนอยู่ในช่วง $[0, 160]$ เซนติเมตร เทอมเซตของคนที่มีความสูงปานกลางจะมีโดเมนอยู่ในช่วง $[150, 170]$ เซนติเมตร และเทอมเซตของคนที่สูงจะมีค่าตั้งแต่ 160 เซนติเมตรขึ้นไป

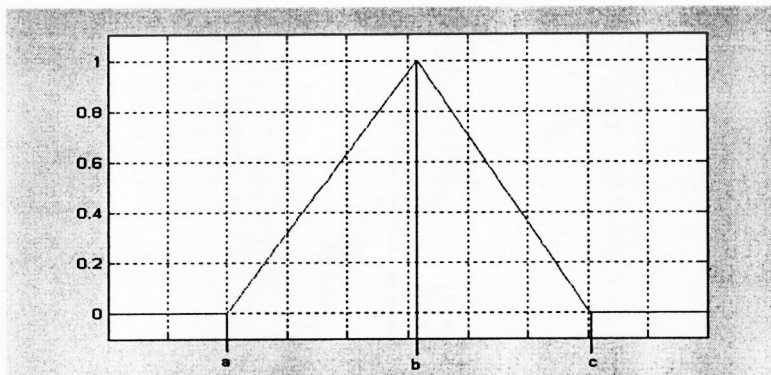


รูปที่ 3.3 กราฟฟังก์ชันความเป็นสมาชิกของตัวแปร “ความสูง” ที่ประกอบไปด้วย 3 เทอมเซต

กราฟฟังก์ชันความเป็นสมาชิกมีหลายแบบที่จะใช้ในการประมาณค่าระดับความเป็นสมาชิกของแต่ละเทอมเซต โดยมีทั้งแบบที่ประมาณค่าเป็นเชิงเส้น เช่น รูปสามเหลี่ยม รูปสี่เหลี่ยมคางหมู และแบบที่ประมาณค่าโดยฟังก์ชันต่อเนื่อง เช่น แบบเส้นโค้งรูปตัว S รูปประฆังคว่ำและเกาส์เซียน เป็นต้น ซึ่งมีรายละเอียดดังนี้

ฟังก์ชันความเป็นสมาชิกแบบสมการเชิงเส้นเป็นฟังก์ชันที่ใช้ในการแปลงค่าความสัมพันธ์ของโดเมนไปยังเรนจ์สำหรับข้อมูลที่มีความสัมพันธ์เป็นเชิงเส้น ดังสมการที่ 3.3

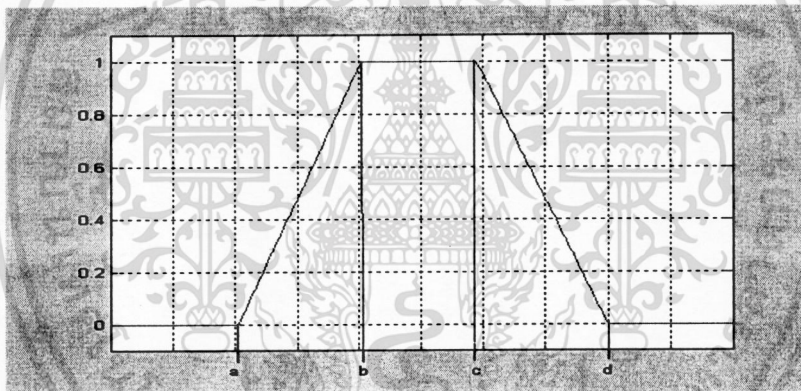
ฟังก์ชันความเป็นสมาชิกแบบรูปสามเหลี่ยม (Triangular membership function) เป็นฟังก์ชันความเป็นสมาชิกที่นิยมใช้กันมากในระบบควบคุมที่ใช้ฟuzzyเซต เทอมเซตที่กำหนดโดยฟังก์ชันนี้จะต้องมีค่าที่เหมาะสมที่สุดอยู่เพียงค่าเดียวที่ทำให้ค่าความเป็นสมาชิกเท่ากับหนึ่ง ส่วนค่าอื่นๆ จะมีค่าความเป็นสมาชิกลดลงเรื่อยๆ เมื่อยิ่งห่างจากค่านี้มากขึ้น ในตัวอย่างของตัวแปรความสูง เทอมเซต “ปานกลาง” มีฟังก์ชันความเป็นสมาชิกเป็นรูปสามเหลี่ยม ถ้ากำหนดให้ $a \leq b \leq c$ เมื่อ a , b , และ c เป็นเลขจำนวนจริงใดๆ ดังนั้นสมการของฟังก์ชันความเป็นสมาชิกแบบรูปสามเหลี่ยมสามารถกำหนดได้ดังสมการที่ 3.4



เอกสารนี้เป็นเอกสาร **รูปที่ 3.4** กราฟแสดงฟังก์ชันความเป็นสมาชิกแบบรูปสามเหลี่ยม ซึ่งประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$$\text{triangle}(x; a, b, c) = \begin{cases} 0 & x \leq a \\ \frac{x-a}{b-a} & a \leq x \leq b \\ \frac{c-x}{c-b} & b \leq x \leq c \\ 0 & c \leq x \end{cases} \quad (3.4)$$

ฟังก์ชันความเป็นสมาชิกแบบรูปสี่เหลี่ยมคางหมู (Trapezoidal membership function) เป็นฟังก์ชันความเป็นสมาชิกที่กำหนดด้วยสมการรูปสี่เหลี่ยมคางหมู นิยมใช้ในระบบควบคุมที่ใช้ฟังก์ชันเชิงเส้นกัน แต่จะพบน้อยกว่าฟังก์ชันความเป็นสมาชิกที่กำหนดด้วยสมการรูปสามเหลี่ยม ฟังก์ชันความเป็นสมาชิกที่กำหนดด้วยสมการรูปสี่เหลี่ยมคางหมูจะต้องมีช่วงของค่าที่เหมาะสมมากที่สุดอยู่ในกลุ่มหนึ่งที่ทำให้ค่าความเป็นสมาชิกของเทอมเซตนั้นๆ นอกจากนั้นจะมีค่าความเป็นสมาชิกน้อยลงเรื่อยๆ เมื่อยิ่งห่างจากข้อมูลกลุ่มนี้ ถ้ากำหนดให้ $a \leq b \leq c \leq d$ เมื่อ $a, b, c,$ และ d เมื่อ $a, b, c,$ และ d เป็นเลขจำนวนจริงใดๆ ดังนั้นสมการของฟังก์ชันความเป็นสมาชิกแบบรูปสี่เหลี่ยมคางหมูสามารถกำหนดได้ดังสมการที่ 3.5

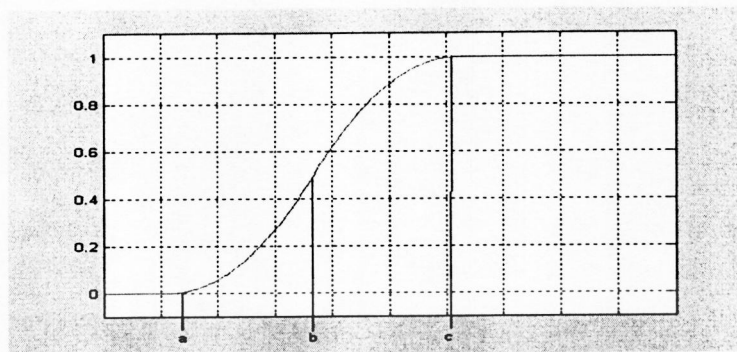


รูปที่ 3.5 กราฟแสดงฟังก์ชันความเป็นสมาชิกแบบรูปสี่เหลี่ยมคางหมู

$$\text{trapezoid}(x; a, b, c, d) = \begin{cases} 0 & x \leq a \\ \frac{x-a}{b-a} & a \leq x \leq b \\ 1 & b \leq x \leq c \\ \frac{d-x}{d-c} & c \leq x \leq d \\ 0 & d \leq x \end{cases} \quad (3.5)$$

ฟังก์ชันความเป็นสมาชิกแบบเส้นโค้งรูปตัว S (Sigmoidal membership function) เป็นฟังก์ชันความเป็นสมาชิกที่กำหนดความสัมพันธ์ระหว่างโดเมนไปยังเรนจ์แบบไม่เป็นเชิงเส้น ตัวแปรที่เหมาะสมที่จะใช้ฟังก์ชันนี้คือพวกอายุการใช้งานของอุปกรณ์ต่างๆ ตัวแปรส่วนใหญ่ที่ใช้ในทางการประมวลผลทางภาพเป็นต้น ซึ่งตัวแปรพวกนี้มักจะมีความสัมพันธ์แบบไม่เป็นเชิงเส้นที่เป็นรูปตัว S หรือ S⁻¹ ถ้ากำหนดให้ $a \leq b \leq c$ เมื่อ $a, b,$ และ c เป็นเลขจำนวนจริงใดๆ ดังนั้น

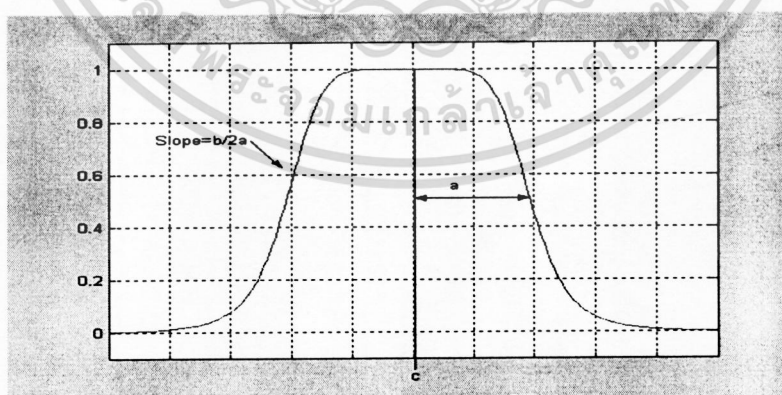
สมการของฟังก์ชันความเป็นสมาชิกของเส้นโค้งรูปตัว S สามารถกำหนดได้ดังสมการที่ 3.6 และ $S' = 1 - \mu_s(x)$ โดยที่ $\mu_s(x)$ ได้จากสมการที่ 3.6



รูปที่ 3.6 กราฟแสดงฟังก์ชันความเป็นสมาชิกแบบรูปตัว S

$$\text{sigmoid}(x; a, b, c) = \begin{cases} 0 & x \leq a \\ 2((x-a)/(c-a))^2 & a \leq x \leq b \\ 1 - 2((x-a)/(c-a))^2 & b \leq x \leq c \\ 1 & x \geq c \end{cases} \quad (3.6)$$

ฟังก์ชันความเป็นสมาชิกแบบรูประฆังคว่ำ (Bell membership function) เป็นฟังก์ชันความเป็นสมาชิกที่มีคุณสมบัติของข้อมูลคล้ายกับฟังก์ชันความเป็นสมาชิกแบบรูปสามเหลี่ยม แต่ข้อมูลของโดเมนที่แปลงไปยังเรนจ์จะมีความสัมพันธ์ในลักษณะไม่เป็นเชิงเส้น ฟังก์ชันความเป็นสมาชิกแบบรูประฆังคว่ำสามารถกำหนดได้ดังสมการที่ 3.7



รูปที่ 3.7 กราฟแสดงฟังก์ชันความเป็นสมาชิกแบบรูประฆังคว่ำ

$$\text{bell}(x; a, b, c) = \begin{cases} 0 & x \leq a \\ 2((x-a)/(c-a))^2 & a < x \leq b \\ 1 - 2((x-a)/(c-a))^2 & b < x \leq c \\ 0 & x > c \end{cases} \quad (3.7)$$

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์ไว้ใช้ในวงจำกัดให้นำไปใช้ประโยชน์ในการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

โดยที่ $b = \frac{c-a}{2}$

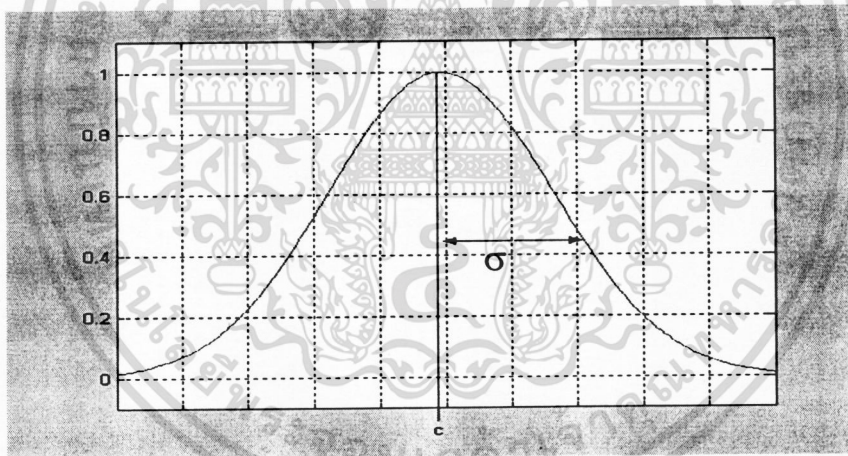
ดังนั้น ฟังก์ชันความเป็นสมาชิกแบบรูปประฆังคว่ำยังสามารถกำหนดได้โดยสมการที่ 3.8

$$\pi(x) = \begin{cases} bell(x; c-b, c-b/2, c) & , x \leq c \\ 1 - bell(x; c-b, c-b/2, c+b) & , x > c \end{cases} \quad (3.8)$$

ฟังก์ชันความเป็นสมาชิกแบบเกาส์เซียน (Gaussian membership function) สามารถกำหนดได้โดยสมการที่ 3.9

$$gaussian(x; c, \sigma) = e^{-\frac{(x-c)^2}{2\sigma^2}} \quad (3.9)$$

โดยที่ c เป็นค่ากลางและ σ เป็นความกว้างของฟังก์ชันความเป็นสมาชิกแบบเกาส์เซียนตามลำดับ



รูปที่ 3.8 กราฟแสดงฟังก์ชันความเป็นสมาชิกแบบรูป Gaussian

3.4 สรุป

ทฤษฎีฟuzzyเซตจะระบุการเป็นสมาชิกโดยอาศัยค่าความเป็นสมาชิกซึ่งอยู่ในช่วง 0 ถึง 1 แตกต่างจากค่าความเป็นสมาชิกของเซตธรรมดาที่เป็นได้เพียง 0 หรือ 1 ทำให้เกิดความยืดหยุ่นของข้อมูล โดยได้มีการพัฒนาฟuzzyเซตไปใช้ในงานด้านต่างๆ มากมาย เช่น ทางด้านระบบควบคุม ทางด้านประมวลผลภาพ เป็นต้น โดยในงานวิจัยนี้ได้นำฟuzzyเซตมาใช้ในการแบ่งช่วง

ลักษณะเด่นที่มีความต่อเนื่องของข้อมูลที่เหมาะสม

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 4

ราฟเซ็ต

ทฤษฎีราฟเซ็ตได้ถูกพัฒนาขึ้นโดย Zdzislaw Pawlak ในปี ค.ศ. 1980 เป็นวิธีการคณิตศาสตร์แบบใหม่ที่ใช้วิเคราะห์ข้อมูลที่มีความคลุมเครือและความไม่แน่นอน ทฤษฎีราฟเซ็ตเกิดมาจากแนวคิดที่จะนำเอาสารสนเทศ (ข้อมูล, ความรู้) ที่เกี่ยวกับทุกๆ วัตถุในระบบที่สนใจมาเป็นตัวอธิบายคุณลักษณะของวัตถุนั้น ตัวอย่างเช่น ถ้าวัตถุในที่นี้คือคนไข้ที่ป่วยด้วยโรคใดโรคหนึ่ง อาการต่างๆ ที่เกิดจากโรคนั้นๆ จะประกอบกันขึ้นเป็นสารสนเทศเกี่ยวกับคนไข้คนนั้น โดยจะกล่าวได้ว่าวัตถุนั้นเหมือนกัน ถ้าวัตถุเหล่านั้นมีสารสนเทศที่บรรยายคุณลักษณะของมันเหมือนกัน ทฤษฎีของราฟเซ็ตได้ถูกพัฒนาไปใช้ในหลายๆ ด้าน เช่น ระบบผู้เชี่ยวชาญในการตัดสินใจ การจำแนกวัตถุ การคัดเลือกคุณลักษณะและการดึงคุณลักษณะ เป็นต้น

4.1 การแสดงค่าความรู้

ในทฤษฎีราฟเซ็ตประกอบด้วยระบบที่เป็นตัวแทนการแสดงความรู้ (Knowledge representation) อยู่ 2 ระบบ ได้แก่ ระบบสารสนเทศ (Information system) และระบบการตัดสินใจ (Decision system) ซึ่งในระบบทั้งสองจะมีเพียงข้อมูลดิบของวัตถุเท่านั้น ไม่มีการแปลหรือตีความ หรือการตั้งสมมติฐานเกี่ยวกับวัตถุนั้นๆ แต่ประการใด โดยจะแสดงข้อมูลดิบอยู่ในรูปของตารางข้อมูล โดยที่แถวจะแทนวัตถุ (objects) และคอลัมน์แทนคุณลักษณะ (attributes) ของวัตถุ ตารางที่ 4.1 จะแสดงรูปแบบของระบบสารสนเทศและระบบการตัดสินใจ

ตารางที่ 4.1 รูปแบบของระบบสารสนเทศและระบบการตัดสินใจ

	a_1	...	a_j	...	a_m
x_1	a_1x_1		a_jx_1		a_mx_1
\vdots		\ddots			
x_i			a_jx_i		a_mx_i
\vdots				\ddots	
x_n					a_mx_n

โดยที่ a_{j,x_i} คือคุณลักษณะที่ j ของวัตถุที่ i โดย $i=1,2,3,\dots,n$ และ $j=1,2,3,\dots,m$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

4.1.1 ระบบสารสนเทศ

ระบบสารสนเทศ คือระบบแสดงความรู้ (knowledge) ในขั้นพื้นฐานที่สุดที่ประกอบด้วยค่าคุณลักษณะ ระบบสารสนเทศ S ประกอบด้วย 4 tuples $S = (U, A, V, f)$ โดย $U = \{x_1, \dots, x_n\}$ คือเซตของวัตถุ (n คือจำนวนของวัตถุ) ซึ่งต่อไปนี้จะเรียกว่า “เซตเอกภพสัมพัทธ์” (universe) $A = \{a_1, \dots, a_m\}$ คือเซตของคุณลักษณะ (m คือจำนวนของคุณลักษณะ) $V = \bigcup_{a \in A} V_a$ และ V_a คือเซตค่าของคุณลักษณะ a ที่อยู่ใน A และ f คือฟังก์ชัน $f: U \times V \rightarrow V_a$

ตารางที่ 4.2 จะแสดงตัวอย่างของระบบสารสนเทศที่ไม่ซับซ้อน ซึ่งระบบสารสนเทศนี้ประกอบด้วยคนไข้หรือวัตถุ 7 คน และคุณลักษณะ 2 ประเภทคือ Age และ LEMS (Lower Extremity Motor Score)

ตารางที่ 4.2 ตัวอย่างของระบบสารสนเทศของคนไข้

	Age	LEMS
x_1	16-30	50
x_2	16-30	0
x_3	31-45	1-25
x_4	31-45	1-25
x_5	46-60	26-49
x_6	16-30	26-49
x_7	46-60	26-49

จากตารางที่ 4.2 จะพบปัญหาว่าไม่สามารถแยกคนไข้ออกจากกันได้ เนื่องจากคนไข้มีอาการเหมือนกัน เช่น ไม่สามารถแยกคนไข้ x_3 และ x_4 กับ x_5 และ x_7 ออกจากกันได้ เนื่องจากคนไข้เหล่านี้มีอาการเหมือนกัน

4.1.2 ระบบการตัดสินใจ

ระบบการตัดสินใจนี้จะมีลักษณะเหมือนกับระบบสารสนเทศ แต่จะแตกต่างกันตรงที่ระบบการตัดสินใจจะมีคุณลักษณะอยู่ 2 ประเภท คือคุณลักษณะที่ใช้ในการสร้างเงื่อนไข (Conditional attributes) และคุณลักษณะที่ใช้ในการตัดสินใจ (Decision attributes) ดังที่ได้กล่าวมาแล้วข้างต้นว่าในระบบสารสนเทศ ข้อมูลจะไม่ถูกต้องความ แต่เมื่อมีการแยกประเภทของวัตถุในระบบโดยผู้เชี่ยวชาญ และมีการกำหนดค่าคุณลักษณะ ก็จะทำให้ระบบสารสนเทศนั้นกลายเป็นระบบการตัดสินใจ

ระบบการตัดสินใจ S ประกอบด้วยคู่ลำดับ $S = (U, A \cup \{d\})$ ซึ่งเป็นระบบสารสนเทศที่แบ่งเอกสารนี้เป็นเอกสารที่ส่งงานไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า คุณลักษณะออกเป็น 2 กลุ่มที่ไม่เกี่ยวเนื่องกัน โดยที่ $d \notin A$ เป็นคุณลักษณะที่ใช้ในการตัดสินใจ ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

(Decision attributes) และ A เป็นคุณลักษณะที่ใช้ในการสร้างเงื่อนไข (Conditional attributes) และ $A \cap \{d\} = \emptyset$

ตารางที่ 4.3 จะแสดงตัวอย่างของระบบการตัดสินใจที่ประกอบด้วยคนไข้ 7 คน เหมือนกับระบบสารสนเทศที่แสดงในตารางที่ 4.2 แต่จะแตกต่างตรงที่มีคุณลักษณะที่ใช้ในการตัดสินใจที่ชื่อ *Walk* เพิ่มขึ้นมาอีกหนึ่งคุณลักษณะ

ตารางที่ 4.3 ตัวอย่างของระบบการตัดสินใจ

	<i>Age</i>	<i>LEMS</i>	<i>Walk</i>
x_1	16-30	50	Yes
x_2	16-30	0	No
x_3	31-45	1-25	No
x_4	31-45	1-25	Yes
x_5	46-60	26-49	No
x_6	16-30	26-49	Yes
x_7	46-60	26-49	No

จากตารางที่ 4.3 จะพบปัญหาว่าคนไข้ x_3 และ x_4 กับ x_5 และ x_7 ยังมีอาการเหมือนกัน แต่คนไข้คู่แรกสามารถแยกออกจากกันได้ เนื่องจากคุณลักษณะที่ใช้ในการตัดสินใจมีค่าไม่เหมือนกัน ส่วนคนไข้คู่ที่สองไม่สามารถแยกออกจากกันได้ เนื่องจากคุณลักษณะที่ใช้ในการตัดสินใจมีค่าเหมือนกัน

4.2 ความคล้ายกันของวัตถุ

ในระบบการตัดสินใจจะมีทั้งวัตถุที่มีคุณลักษณะที่แตกต่าง และวัตถุที่มีคุณลักษณะที่เหมือนกัน โดยการระบุความแตกต่างของวัตถุเหล่านั้นจะยึดเอาค่าคุณลักษณะของวัตถุเป็นหลัก หรือเป็นข้อมูลในการจำแนก

ความเหมือนกันระหว่างวัตถุในระบบอันเนื่องมาจากค่าคุณลักษณะที่ถูกนำมาใช้งานนั้น (โดยที่ในความเป็นจริงวัตถุอาจจะไม่มีความเหมือนกันเลย) จะถูกเรียกว่า “ความสัมพันธ์แบบคล้ายกัน” (Indiscernibility relation) ซึ่งเป็นความสัมพันธ์ที่จะทำให้สามารถแบ่งเซตเอกภพสัมพัทธ์ออกเป็นซับเซตย่อยๆ ที่ไม่เกี่ยวเนื่องกันได้ โดยในแต่ละซับเซตนั้นจะมีสมาชิกเป็นวัตถุที่มีความเหมือนกันหรือมีความสัมพันธ์แบบคล้ายกัน โดยวัตถุจากค่าคุณลักษณะที่ถูกเลือกมาใช้

ความสัมพันธ์แบบคล้ายกันนี้สามารถนิยามได้ดังนี้ ให้ $S = (U, A)$ เป็นระบบสารสนเทศ เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ทุกๆ ซับเซตของค่าคุณลักษณะ $B \subseteq A$ จะกำหนดความสัมพันธ์เท่าเทียม (Equivalent relation) ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$IND_S(B)$ ซึ่งต่อไปจะเรียกว่า “ความสัมพันธ์แบบคล้ายกันที่กำหนดด้วยเซต B ” (B-indiscernibility relation) ซึ่งสามารถกำหนดได้ดังสมการที่ 4.1

$$IND_S(B) = \{ (x, x') \in U^2 : \text{for every } a \in B, a(x) = a(x') \} \quad (4.1)$$

โดยที่ a คือค่าคุณลักษณะ และ $a(x)$ คือค่าของวัตถุ x ณ คุณลักษณะ a

แนวคิดของความสัมพันธ์แบบคล้ายกัน คือการเลือกเซตของค่าคุณลักษณะ $B \subseteq A$ ที่ทำให้เกิดการแบ่งแยกเซตเอกภพสัมพัทธ์ออกเป็นเซตย่อยอย่างสมบูรณ์ โดยที่สมาชิกของแต่ละเซตย่อยนั้นจะไม่สามารถแบ่งแยกได้อีกโดยอาศัยเพียงค่าคุณลักษณะที่อยู่ในเซต B และแต่ละเซตย่อยที่เกิดขึ้นจากการแบ่งเซตเอกภพสัมพัทธ์จะถูกเรียกว่า “เซตของวัตถุที่คล้ายกัน” (Equivalent classes) เซตของวัตถุที่คล้ายกันที่กำหนดด้วยเซต B จะถูกเขียนแทนด้วยสัญลักษณ์ $[x]_B$ และนิยามได้ดังนี้

ให้ $S = (U, A)$ เป็นระบบสารสนเทศ และ $B \subseteq A$ จากนั้นจะนิยามเซตของวัตถุที่คล้ายกันที่กำหนดด้วยเซต B ได้ดังสมการที่ 4.2

$$[x]_B = \{ y \in U \mid (x, y) \in IND_S(B) \} \quad (4.2)$$

โดยที่ x และ y คือสมาชิกของเซตเอกภพสัมพัทธ์

ต่อไปจะแสดงการกำหนดความสัมพันธ์แบบคล้ายกันในระบบการตัดสินใจ จากตารางที่ 4.3 กำหนดให้ $\{Age\}$, $\{LEMS\}$, และ $\{Age, LEMS\}$ เป็นเซตของคุณลักษณะที่ใช้ในการสร้างเงื่อนไข ซึ่งถ้าพิจารณาเซต $\{LEMS\}$ คนไข้ x_3 และ x_4 จะอยู่ในเซตของวัตถุที่คล้ายกันที่กำหนดด้วยคุณลักษณะ $LEMS$ และมีความสัมพันธ์แบบคล้ายกัน และคนไข้ x_5, x_6 , และ x_7 จะอยู่ในเซตของวัตถุที่คล้ายกันเซตอื่น ความสัมพันธ์ IND จะกำหนดพาร์ติชันของเซตเอกภพสัมพัทธ์ได้ 3 แบบดังนี้

$$IND(\{Age\}) = \{ \{x_1, x_2, x_6\}, \{x_3, x_4\}, \{x_5, x_7\} \}$$

$$IND(\{LEMS\}) = \{ \{x_1\}, \{x_2\}, \{x_3, x_4\}, \{x_5, x_6, x_7\} \}$$

$$IND(\{Age, LEMS\}) = \{ \{x_1\}, \{x_2\}, \{x_3, x_4\}, \{x_5, x_7\}, \{x_6\} \}$$

4.3 การประมาณค่าของเซต

ในหัวข้อนี้จะกล่าวถึงทฤษฎีของราฟเซตที่ใช้ดำเนินการกับข้อมูลหรือแนวคิดที่มีความ
เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์สำหรับใช้ในเพื่อการศึกษาเท่านั้น เมื่ออนุญาตให้ไปเผยแพร่บนเว็บไซต์
กลุ่มเครือข่าย
ไม่มีส่วนไหนใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ความสัมพันธ์แบบคล้ายกัน (equivalence relation) ทำให้เกิดพาร์ติชันของเซตเอกภพสัมพัทธ์ที่สามารถนำไปใช้สร้างซับเซตของเซตเอกภพสัมพัทธ์ขึ้นมาใหม่ได้ โดยจะให้ความสนใจเฉพาะซับเซตที่มีค่าคุณลักษณะที่เหมือนกันหรืออยู่ในกลุ่มเดียวกัน แต่อย่างไรก็ตาม แนวคิดอย่างเช่น “Walk” ไม่สามารถกำหนดโดยใช้คุณลักษณะที่มีอยู่ได้ง่ายๆ เช่น เซตของคนไข้ที่มีผลที่ชัดเจน (Positive outcome) ไม่สามารถกำหนดโดยใช้คุณลักษณะที่มีในตารางที่ 4.3 ได้ง่าย โดยคนไข้ที่มีปัญหาคือ คนไข้ x_3 และ x_4 อธิบายได้อีกอย่างว่า เป็นไม่ได้ที่จะอธิบายลักษณะที่ถูกต้องของคนไข้จากตารางได้ จากปัญหานี้จึงได้เกิดแนวคิดของกราฟเซตขึ้น ซึ่งจะทำให้มีความเป็นไปได้ที่จะจำแนกคนไข้ที่มีผลที่ชัดเจนและไม่ชัดเจน และสุดท้ายคนไข้ที่อยู่ในขอบเขต (boundary) ระหว่างอาการที่ชัดเจนได้อย่างแน่นอน แม้ว่าจะไม่สามารถจำแนกคนไข้ได้ง่ายก็ตาม แนวคิดเหล่านี้จะถูกแสดงดังต่อไปนี้

ให้ $S = (U, A)$ เป็นระบบสารสนเทศ $B \subseteq A$ เป็นเซตของค่าคุณลักษณะที่ถูกเลือก และ $X \subseteq U$ เป็นเซตของวัตถุ เราสามารถประมาณค่าของเซต X โดยใช้เพียงสารสนเทศที่มีในเซต B ด้วยการสร้าง B -Lower Approximation และ B -Upper Approximation ของ X เขียนแทนด้วย \underline{BX} และ \overline{BX} ตามลำดับ ซึ่งถูกนิยามได้โดยสมการต่อไปนี้

$$\underline{BX} = \{x \in U : [x]_B \subseteq X\} \quad (4.3)$$

$$\overline{BX} = \{x \in U : [x]_B \cap X \neq \emptyset\} \quad (4.4)$$

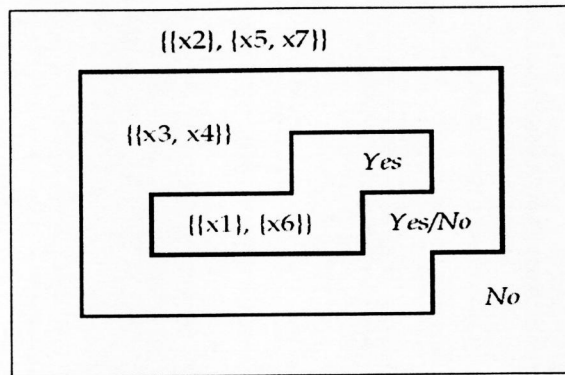
จากสมการที่ 4.3 และ 4.4 วัตถุที่เป็นสมาชิกในเซต \underline{BX} สามารถถูกจำแนกโดยใช้ความรู้ที่มีอยู่ในเซต B ว่าเป็นสมาชิกของเซต X อย่างแน่นอน ส่วนวัตถุที่เป็นสมาชิกในเซต \overline{BX} สามารถถูกจำแนกโดยใช้ความรู้ที่มีอยู่ในเซต B ว่าอาจจะเป็นไปได้ที่จะเป็นสมาชิกของเซต X จากที่กล่าวมาจะได้เซต B -Boundary Region ของ X เขียนแทนด้วย $BN_B(X)$ ซึ่งถูกนิยามได้โดยสมการต่อไปนี้

$$BN_B(X) = \overline{BX} - \underline{BX} \quad (4.5)$$

จากสมการที่ 4.5 วัตถุที่เป็นสมาชิกในเซต $BN_B(X)$ ไม่สามารถจำแนกโดยใช้ความรู้ในเซต B ได้ว่าเป็นสมาชิกของเซต X อย่างแน่นอน

จากตารางที่ 4.3 อาการทั่วไปส่วนมากจะถูกวิเคราะห์หาค่าความหมายของผลหรือกลุ่มของการตัดสินใจในเทอมของคุณลักษณะที่ใช้ในการสร้างเงื่อนไข เช่น กำหนดให้ $W = \{x | \text{Walk}(x) = \text{Yes}\}$ จากนั้นจะได้ค่าประมาณของเซต $\underline{AW} = \{x_1, x_6\}$, $\overline{AW} = \{x_1, x_3, x_4, x_6\}$, และ $BN_A(W) = \{x_3, x_4\}$ ดังนั้นก็สรุปได้ว่าผลของ Walk เป็นกราฟเซต เนื่องจาก boundary region ไม่ได้เป็นเซตว่าง ซึ่งแสดงดังรูปที่ 4.1

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.1 การประมาณค่าเซตคนไข้ที่เดินได้ โดยใช้ Age และ LEMS เป็นคุณลักษณะที่ใช้ในการสร้างเงื่อนไข กลุ่มของคนไข้ที่มีอาการคล้ายกันจะอยู่ใน regions เดียวกัน

นอกจากนี้ ยังสามารถถูกอธิบายลักษณะกราฟเซต โดยใช้ค่าสัมประสิทธิ์ดังสมการต่อไปนี้

$$\alpha_B = \frac{|B(X)|}{|B(X)|} \quad (4.6)$$

ซึ่งเรียกว่า ค่าประมาณของความถูกต้อง (Accuracy of approximation) โดย $|X|$ คือจำนวนสมาชิกของ X ค่าประมาณของความถูกต้องมีค่าอยู่ระหว่าง 0 ถึง 1 ($0 \leq \alpha_B(X) \leq 1$) ถ้า $\alpha_B(X) = 1$ แล้ว X จะเป็น crisp เมื่ออาศัยเซต B (X มีความถูกต้องเมื่อพิจารณาจากเซต B) และถ้า $\alpha_B(X) < 1$ แล้ว X จะเป็น rough เมื่ออาศัยเซต B (X มีความคลุมเครือเมื่อพิจารณาจากเซต B)

4.4 เซตเซตที่มีคุณลักษณะน้อยที่สุดที่ยังคงสามารถจำแนกวัตถุได้

ในหัวข้อก่อนหน้านี ได้กล่าวถึงระบบการตัดสินใจในมิติที่มีข้อมูลน้อยซึ่งจะแสดงกลุ่มของวัตถุที่คล้ายกันโดยใช้คุณลักษณะเท่าที่หาได้ การลดขนาดข้อมูลจะทำได้ตั้งแต่ใช้สมาชิกของกลุ่มวัตถุที่คล้ายกันเพียงกลุ่มเดียวแทนกลุ่มข้อมูลทั้งหมด การลดขนาดข้อมูลอีกอย่างหนึ่งคือการเก็บเพียงคุณลักษณะที่ยังคงความสัมพันธ์แบบคล้ายกันและการประมาณค่าของเซตเหมือนเดิม คุณลักษณะที่ถูกกำจัดจะเป็นคุณลักษณะที่เกินหรือไม่จำเป็น เนื่องจากคุณลักษณะที่ถูกกำจัดไม่ทำให้ประสิทธิภาพการจำแนกต่ำลง โดยทั่วไปจะมีเซตของคุณลักษณะที่ยังคงความสัมพันธ์แบบคล้ายกันและการประมาณค่าของเซตมากมายและเซตที่มีคุณลักษณะน้อยที่สุดจะถูกเรียกว่า reduct

จากตารางที่ 4.4 กำหนดระบบสารสนเทศ $S' = (U, \{Diploma, Experience, French, Reference\} \cup \{Decision\})$ เมื่อพิจารณาเฉพาะคุณลักษณะที่ใช้ในการสร้างเงื่อนไข นั่นคือระบบสารสนเทศ

$$S = (U, \{Diploma, Experience, French, Reference\}) \quad (4.7)$$

แต่ละกลุ่มของวัตถุที่คล้ายกันที่ไม่ซับซ้อนจะมีสมาชิกหนึ่งตัว และเซตที่มีคุณลักษณะน้อยที่สุดคือ $\{Experience, Reference\}$ ซึ่งเป็นเซตที่จำแนกวัตถุด้วยวิธีเดียวกันกับเซตที่มีคุณลักษณะทั้งหมด จากการตรวจสอบความสัมพันธ์แบบคล้ายกันโดยใช้เซตที่มีคุณลักษณะทั้งหมดและเซต $\{Experience, Reference\}$ พบว่าเหมือนกัน ต่อไปจะแสดงการสร้างเซตที่มีคุณลักษณะน้อยที่สุดที่ยังคงความสัมพันธ์แบบคล้ายกันและการประมาณค่าของเซต

ตารางที่ 4.4 ตัวอย่างของระบบตัดสินใจที่ยังไม่ได้ลดคุณลักษณะที่มีชื่อว่า *Hiring*

	<i>Diploma</i>	<i>Experience</i>	<i>French</i>	<i>Reference</i>	<i>Decision</i>
x_1	MBA	Medium	Yes	Excellent	Accept
x_2	MBA	Low	Yes	Neutral	Reject
x_3	MCE	Low	Yes	Good	Reject
x_4	MSc	High	Yes	Neutral	Accept
x_5	MSc	Medium	Yes	Neutral	Reject
x_6	MSc	High	Yes	Excellent	Accept
x_7	MBA	High	No	Good	Accept
x_8	MCE	Low	No	Excellent	Reject

ให้ $S = (U, A)$ เป็นระบบสารสนเทศ จากนั้นจะนิยามแนวคิดที่กล่าวข้างต้นมาได้ดังนี้ reduct ของ S คือเซตที่มีคุณลักษณะน้อยที่สุด $B \subseteq A$ ดังนั้น $IND_S(B) = IND_S(A)$ ซึ่งอธิบายได้อีกอย่างว่า reduct คือเซตที่ประกอบด้วยคุณลักษณะที่มาจากเซต A จำนวนน้อยที่สุดที่ยังคงพาร์ติชันของเซตเอกภพสัมพัทธ์และสามารถจำแนกสมาชิกในเซตเอกภพสัมพัทธ์ได้เหมือนกับเซต A

ให้ S เป็นระบบสารสนเทศที่ประกอบด้วยวัตถุ n ตัว เมทริกซ์ของวัตถุที่คล้ายกัน (Discernibility matrix) ของ S คือ symmetric matrix ขนาด $n \times n$ ที่ประกอบด้วยสมาชิก c_{ij} ซึ่งสามารถนิยามได้ดังสมการต่อไปนี้

$$c_{ij} = \begin{cases} \{a \in A : a(x_i) \neq a(x_j)\}, & D(x_i) \neq D(x_j) \\ \emptyset, & D(x_i) = D(x_j) \end{cases} \quad (4.8)$$

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามให้ตัดแปลงเนื้อหา และต้องอ้างอิงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

โดยที่ $a(x)$ คือค่าของวัตถุ x ณ. คุณลักษณะ a และ $D(x)$ คือค่าของ x ณ. เซ็ตของคุณลักษณะที่ใช้ในการตัดสินใจ D

จากสมการที่ 4.8 c_{ij} ซึ่งเป็นสมาชิกที่อยู่ในตำแหน่งแถวที่ i และหลักที่ j ของ Discernibility matrix คือเซตที่ประกอบด้วยคุณลักษณะที่มีค่าไม่เท่ากันของวัตถุ x_i และ x_j ถ้าค่าตัดสินใจของวัตถุ x_i และ x_j มีค่าไม่เท่ากัน และจะเป็นเซตว่าง ถ้าค่าตัดสินใจของวัตถุ x_i และ x_j มีค่าเท่ากัน ตารางที่ 4.5 แสดง Discernibility matrix ที่สร้างจากระบบตัดสินใจ *Hiring* ที่อยู่ในตารางที่ 4.4 โดยใช้อักษรตัวแรกเขียนแทนแต่ละคุณลักษณะ

ตารางที่ 4.5 Discernibility matrix ของระบบตัดสินใจ *Hiring*

$x \in U$	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8
x_1	ϕ							
x_2	e,r	ϕ						
x_3	d,e,r	ϕ	ϕ					
x_4	ϕ	d,e	d,e,r	ϕ				
x_5	d,r	ϕ	ϕ	e	ϕ			
x_6	ϕ	d,e,r	d,e,r	ϕ	e,r	ϕ		
x_7	ϕ	e,f,r	d,e,f	ϕ	d,e,f,r	ϕ	ϕ	
x_8	d,e,f	ϕ	ϕ	d,e,f,r	ϕ	d,e,f	d,e,r	ϕ

Discernibility function f_S ของระบบสารสนเทศ S คือฟังก์ชันบูลีนที่มีตัวแปรจำนวน m ตัว ได้แก่ a_1^*, \dots, a_m^* ที่ตรงกับคุณลักษณะ a_1, \dots, a_m ซึ่งสามารถนิยามได้ดังสมการต่อไปนี้

$$f_S(a_1^*, \dots, a_m^*) = \bigwedge \{ \bigvee c_{ij}^* \mid 1 \leq j < i \leq n, c_{ij} \neq \phi \} \quad (4.9)$$

โดย $c_{ij}^* = \{ a^* \mid a \in c_{ij} \}$ และเซต prime implicants ทั้งหมดของ f_S จะกำหนดเซต reducts ทั้งหมดของระบบสารสนเทศ S

จากสมการที่ 4.9 Discernibility function จะถูกเขียนอยู่ในรูปแบบ CNF (Conjunctive normal form) โดย a เป็นสมาชิกของเซต c_{ij} และ c_{ij} ไม่เป็นเซตว่าง การหา Discernibility function ทำได้โดยนำสมาชิก c_{ij} ของ Discernibility matrix มาเขียนอยู่ในรูปแบบ DNF (Disjunctive normal form) ที่จากนั้นนำแต่ละเทอมมาเขียนรวมกันในรูปแบบ CNF เมื่อทำการลดรูปโดยใช้

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ทฤษฎี absorption law และ distribution law จะได้ Discernibility function ที่ประกอบด้วยเซตของ reducts ทั้งหมดของระบบสารสนเทศ S

จาก Discernibility matrix ของระบบตัดสินใจ *Hiring* ที่แสดงในตารางที่ 4.5 สามารถเขียน Discernibility function f_S ได้ดังนี้

$$\begin{aligned} f_S(d,e,f,r) = & (e \vee r) \wedge (d \vee e \vee r) \wedge (d \vee r) \wedge (d \vee e \vee f) \wedge (d \vee e) \wedge \\ & (d \vee e \vee r) \wedge (e \vee f \vee r) \wedge (d \vee e \vee r) \wedge (d \vee e \vee r) \wedge \\ & (d \vee e \vee f) \wedge (e) \wedge (d \vee e \vee f \vee r) \wedge (e \vee r) \wedge \\ & (d \vee e \vee f \vee r) \wedge (d \vee e \vee f) \wedge (d \vee e \vee r) \end{aligned}$$

โดยที่แต่ละเทอมที่อยู่ในวงเล็บคือ conjunction ที่อยู่ในนิพจน์บูลีน และตัวแปรบูลีนที่แทนด้วยตัวอักษรหนึ่งตัวจะตรงกับชื่อของคุณลักษณะ หลังจากทำการลดรูปจะได้ฟังก์ชัน $f_S(d,e,f,r) = ed \vee er$ (เครื่องหมาย ed และ er จะแทน $e \wedge d$ และ $e \wedge r$ ตามลำดับ) และจากฟังก์ชันนี้สรุปได้ว่าระบบตัดสินใจ *Hiring* ที่แสดงในตารางที่ 4.4 จะมี reduct หรือเซตที่มีคุณลักษณะน้อยที่สุดที่สามารถจำแนกสมาชิกในเซตเอกภพสัมพัทธ์ได้เหมือนกับเซต A สองเซต คือ $\{Experience, Reference\}$ และ $\{Diploma, Experience\}$

4.5 การขึ้นต่อกันของคุณลักษณะ

สิ่งที่สำคัญอีกอย่างหนึ่งในการวิเคราะห์ข้อมูล คือการคำนวณหาการขึ้นต่อกันระหว่างคุณลักษณะ ซึ่งจะเห็นได้อย่างชัดเจนว่าเซตของคุณลักษณะ D จะขึ้นกับเซตของคุณลักษณะ C โดยทั้งหมด เขียนแทนด้วย $C \Rightarrow D$ ถ้าค่าทั้งหมดของคุณลักษณะที่มาจาก D ถูกหาได้โดยค่าของคุณลักษณะที่มาจาก C อธิบายได้อีกอย่างว่า D จะขึ้นกับ C โดยทั้งหมด ถ้ายังมีการขึ้นต่อกันระหว่างของ D และ C

การขึ้นต่อกันของคุณลักษณะสามารถนิยามได้ด้วยวิธีต่อไปนี้ ให้ D และ C เป็นซับเซตของ A โดยจะสรุปว่าซับเซต D ขึ้นอยู่กับซับเซต C ในระดับ k ($0 \leq k \leq 1$) เขียนแทนด้วย $C \Rightarrow_k D$ ถ้า

$$k = \frac{|POS_C(D)|}{|U|} \quad (4.10)$$

โดย

$$POS_C(D) = \bigcup_{x \in U/D} C(x) \quad (4.11)$$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เรียกว่า “*positive region*” ของพาร์ติชัน UID ซึ่งสอดคล้องกับชั้นเซต C นั่นคือเซตของสมาชิกของเซตเอกภพสัมพัทธ์ทั้งหมดที่สามารถถูกจำแนกเป็นบล็อกโคบล็อกหนึ่งของพาร์ติชัน UID โดยอาศัยชั้นเซต C เห็นได้ชัดว่า

$$k = \sum_{x \in UID} \frac{|C(x)|}{|U|} \quad (4.12)$$

ถ้า $k=1$ จะสรุปได้ว่าชั้นเซต D ขึ้นอยู่กับชั้นเซต C โดยทั้งหมด และถ้า $k < 1$ จะสรุปได้ว่าชั้นเซต D ขึ้นอยู่กับชั้นเซต C เป็นบางส่วน (ในระดับ k)

ค่าสัมประสิทธิ์ k จะแสดงอัตราส่วนของสมาชิกทั้งหมดของเซตเอกภพสัมพัทธ์ที่สามารถจำแนกบล็อกของพาร์ติชัน UID ได้ถูกต้องโดยการใช้คุณลักษณะ C และเรียกว่าระดับการขึ้นอยู่กับคุณลักษณะ (degree of the dependency)

4.6 กฎการตัดสินใจ

จากหัวข้อที่ผ่านมาทำให้เข้าใจได้ว่า *reducts* สามารถนำไปใช้สร้างกฎการตัดสินใจที่มีจำนวนน้อยที่สุดได้ กฎจะถูกสร้างโดยการวาง *reducts* ซ้อนบนตารางตัดสินใจและทำนายค่าได้หลังจากคำนวณ *reducts* ได้แล้ว

จากตารางที่ 4.4 เมื่อให้ $\{Diploma, Experience\}$ เป็น *reduct* จากนั้นกฎที่ทำนายว่าเป็นผู้สมัครคนแรกคือ “if *Diploma* is *MBA* and *Experience* is *Medium* then *Decision* in *Accept*”

ให้ $S=(U, A \cup \{d\})$ เป็นระบบตัดสินใจ และ $V = \bigcup \{V_a : a \in A\} \cup V_d$ กฎที่สั้นที่สุด (atomic formulae) ที่ครอบคลุม $B \subseteq A \cup \{d\}$ และ V จะถูกเขียนอยู่ในรูปแบบ $a = v$ และเรียกว่า *descriptors* ที่ครอบคลุม B และ V โดยที่ $a \in B$ และ $v \in V_a$ เซตของกฎที่ครอบคลุม B และ V คือเซตที่มีสมาชิกน้อยที่สุดที่ประกอบด้วยกฎที่สั้นที่สุดที่ครอบคลุม B และ V ทั้งหมด และถูกนำมาเขียนรวมกันโดย propositional connectives \wedge (conjunction) \vee (disjunction) และ \neg (negation)

ให้ $l \in F(B, V)$ ที่เขียนอยู่ในรูปแบบ $(a_1=v_1) \wedge \dots \wedge (a_l=v_l)$ โดย $v_i \in V_{a_i}$ (เมื่อ $i=1, \dots, l$) และ $\{a_1, \dots, a_l\} \subseteq A$ จากนั้น $\|l_s\|$ ถูกใช้แสดงความหมายของ l ในตารางการตัดสินใจ S คือเซตของวัตถุทั้งหมดที่อยู่ใน U ตามคุณสมบัติ l และถูกนิยามได้ดังต่อไปนี้

1. ถ้า l ถูกเขียนอยู่ในรูปแบบ $a = v$ จากนั้น $\|l_s\| = \{x \in U \mid a(x)=v\}$
2. $\|l_s \wedge l'_s\| = \|l_s\| \cap \|l'_s\|$; $\|l_s \vee l'_s\| = \|l_s\| \cup \|l'_s\|$; $\|\neg l_s\| = U - \|l_s\|$

เซต $F(B, V)$ ถูกเรียกว่า เซตที่ประกอบด้วย *conditional formulae* ของ S และถูกเขียนแทนด้วย

เอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

กฎการตัดสินใจของระบบการตัดสินใจ S จะถูกเขียนอยู่ในรูปแบบ $l \Rightarrow d = v$ โดย $l \in C(B, V)$, $v \in V_d$, และ $\|l_s\| \neq 0$ กฎ l และ $d = v$ หมายถึง predecessor และ successor ของกฎการตัดสินใจ $l \Rightarrow d = v$ ตามลำดับ

กฎการตัดสินใจ $l \Rightarrow d = v$ จะเป็นจริง (true) ถ้า $\|l_s\| \subseteq \|d = v_s\|$ โดยที่ $\|l_s\|$ คือเซตของวัตถุที่ตรง (matching) กับกฎการตัดสินใจ และ $\|l_s\| \cap \|d = v_s\|$ คือเซตของวัตถุที่สนับสนุนหรือยอมรับ (supporting) กฎ ตัวอย่างเช่น กฎบางส่วนที่ได้จากการวิเคราะห์ตารางที่ 4.4 มีลักษณะดังนี้

$Diploma = MBA \quad \wedge \quad Experience = Medium \quad \Rightarrow \quad Decision = Accept$

$Experience = Low \quad \wedge \quad Reference = Good \quad \Rightarrow \quad Decision = Reject$

$Diploma = MSc \quad \wedge \quad Experience = Medium \quad \Rightarrow \quad Decision = Accept$

จากตารางที่ 4.4 สรุปได้ว่ากฎสองข้อแรกเป็นจริง ส่วนกฎข้อที่สามไม่เป็นจริง

การสร้างกฎจะใช้แฟกเตอร์หลายตัว ตัวอย่างเช่น support ของกฎการตัดสินใจ คือจำนวนของวัตถุที่ตรงกับ predecessor ของกฎ ค่าตัวเลขที่เกี่ยวข้องกับความถี่หลายค่าจะถูกคำนวณโดยอาศัยการนับอย่างเช่นค่าสัมประสิทธิ์ความถูกต้อง (Accuracy coefficient) จะมีค่าเท่ากับสมการต่อไปนี้

$$acc = \frac{\|l_s\| \cap \|d = v_s\|}{\|l_s\|} \quad (4.13)$$

เมื่อ acc คือค่าสัมประสิทธิ์ความถูกต้อง

4.7 สรุป

ทฤษฎีกราฟเซตเป็นวิธีทางคณิตศาสตร์ที่สามารถวิเคราะห์ความคลุมเครือของข้อมูล การหาขอบเขตบนและขอบเขตล่างของกลุ่มข้อมูล เพื่อนำมาคำนวณค่าการขึ้นต่อกันของข้อมูล และจากค่าของการขึ้นต่อกันนี้ เราสามารถนำไปวิเคราะห์หาคุณลักษณะที่ไม่มีความจำเป็นต่อการจำแนกกลุ่มได้ ได้มีการพัฒนางานวิจัยทางด้านกราฟเซตไปใช้ในหลายสาขา เช่น ทางการแพทย์และเภสัชศาสตร์ได้นำกราฟเซตไปใช้ในการวินิจฉัยโรค ทางการเงินและการธนาคารได้นำกราฟเซตไปวิเคราะห์ความเสี่ยงทางการเงินหรือการลงทุน หรือแม้แต่ทางด้านวิศวกรรมก็ได้มีการนำกราฟเซตไปใช้ในระบบควบคุมต่างๆ ซึ่งงานวิจัยนี้ได้เสนอการนำกราฟเซตมาใช้ในการจำแนกการบุกรุกใน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ระบบเครือข่ายคอมพิวเตอร์
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

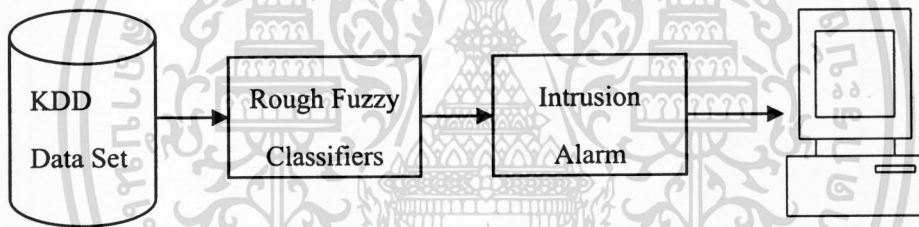
บทที่ 5

การออกแบบระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ ราฟฟี่ซซี่

ในบทนี้จะกล่าวถึงหลักการและวิธีการดำเนินงานวิจัย ของงานวิจัยที่ผู้วิจัยนำเสนอ ซึ่งในส่วนนี้จะนำเสนอถึงวิธีการในการดำเนินงานวิจัย การเตรียมข้อมูลที่ใช้ในการทดลอง แสดงโมเดลการทดลอง และการวัดประสิทธิภาพของระบบ เป็นต้น

5.1 ขั้นตอนการทำงานของระบบ

ในหัวข้อนี้นำเสนอขั้นตอนการทำงานของระบบตรวจจับการบุกรุกที่ผู้วิจัยได้นำเสนอ ซึ่งเป็นการประยุกต์ใช้วิธีราฟฟี่ซซี่ในการจำแนกประเภทการบุกรุกของข้อมูล ดังรูปที่ 5.1



รูปที่ 5.1 แสดงขั้นตอนการทำงานของระบบ

จากรูปที่ 5.1 ส่วนแรก คือข้อมูลที่ใช้ในงานวิจัยนี้เป็นข้อมูลที่ได้จาก KDD Cup 1999 [9] ซึ่งเป็นข้อมูลที่จำลองพฤติกรรมกรบุกรุกที่นิยมใช้ในงานวิจัยที่เกี่ยวกับการตรวจจับผู้บุกรุกจะแสดงรายละเอียดในหัวข้อถัดไป ต่อมาในส่วนที่สองจะเป็นส่วนของการประยุกต์ใช้วิธีราฟฟี่ซซี่ในการสร้างกฎสำหรับจำแนกพฤติกรรมบุกรุกและพฤติกรรมปกติบนเครือข่ายคอมพิวเตอร์ จะแสดงรายละเอียดในหัวข้อ 5.3 ส่วนที่สามจะเป็นส่วนของการแจ้งเตือน ไปสู่เครื่องคอมพิวเตอร์ และในที่สุดท้ายจะเป็นส่วนที่ผู้ดูแลระบบจะได้รับข้อมูลจากการจำแนกพฤติกรรมบุกรุกจากระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟี่ซซี่เพื่อเก็บไว้เป็นข้อมูลในเครื่องคอมพิวเตอร์

5.2 ข้อมูล KDD Cup 1999

ข้อมูล KDD Cup 1999 [9] เป็นข้อมูลที่ประยุกต์มาจากข้อมูลพฤติกรรมกรบุกรุกของเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาค้นคว้า ไม่นับค่าตอบแทนไปประโยชน์ด้านการค้า DARPA 98 ที่ประกอบด้วยทั้งข้อมูลของพฤติกรรมกรบุกรุกและพฤติกรรมปกติ โดยข้อมูลนี้ได้ไม่จำกัดแต่ทุกสิ่งทุกอย่างที่ห้ามมิให้เปิดเผยเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากการวิเคราะห์ระบบเครือข่ายของมหาวิทยาลัยโคลัมเบีย (Columbia University)

5.2.1 ลักษณะข้อมูลของ KDD Cup 1999

ข้อมูลของ KDD Cup 1999 เป็นข้อมูลของการติดต่อสื่อสารบนโปรโตคอล TCP ในระบบเครือข่ายคอมพิวเตอร์ ซึ่งประกอบด้วยพฤติกรรมปกติ และพฤติกรรมการบุกรุก โดยพฤติกรรมการบุกรุกได้แบ่งออกเป็น 4 กลุ่มใหญ่ๆ ดังต่อไปนี้

1. Denial of Service Attacks คือการบุกรุกที่ผู้บุกรุกพยายามทำให้ระบบคอมพิวเตอร์ไม่สามารถให้บริการต่างๆ ได้ เช่น Apache2, Mailbomb, SYN Flood, Ping of death, Smurf, และ UDP Storm เป็นต้น
2. Remote to User Attacks คือการบุกรุกที่ผู้บุกรุกไม่มียูสเซอร์อยู่ในระบบแต่พยายามทำให้ตัวเองเข้าสู่ระบบให้ได้ เช่น Dictionary, Ftp-write, Fuest, Imap, Named, และ Phf sendmail เป็นต้น
3. User to Root Attacks คือการบุกรุกที่ผู้บุกรุกพยายามเข้าสู่ระบบโดยใช้สิทธิ์เท่าเทียมกับ root เช่น Eject, Ffbconfig, Fdformat, และ Loadmodule เป็นต้น
4. Probing Attacks คือการบุกรุกที่ผู้บุกรุกพยายามตรวจหาจุดอ่อนของระบบ เช่น IPsweep, Mscan, Nmap, Saint, และ Satan เป็นต้น

และข้อมูลของ KDD Cup 1999 ยังประกอบด้วยคุณลักษณะที่ได้จากการเชื่อมต่อระบบเครือข่ายจำนวน 41 คุณลักษณะ โดยสามารถแยกคุณลักษณะออกเป็นกลุ่มย่อยได้ 4 กลุ่มดังต่อไปนี้

1. Basic features เป็นคุณลักษณะพื้นฐานที่ได้จากแพคเกจข้อมูลที่สื่อสารในเครือข่าย เช่น เวลาในการเชื่อมต่อ ชนิดของโปรโตคอล ชนิดของการให้บริการ และสถานะแฟล็ก เป็นต้น ประกอบด้วย 9 คุณลักษณะดังนี้

ตารางที่ 5.1 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
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 เป็นคุณลักษณะที่เก็บรวบรวมข้อมูลที่แสดงให้เห็นพฤติกรรมน่าสงสัย เช่น ความผิดพลาดในการล็อกอิน หรือการใช้คำสั่ง “su” เป็นต้น ประกอบด้วย 13 คุณลักษณะดังนี้

ตารางที่ 5.2 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 ประกอบด้วย 9 คุณลักษณะดังนี้

ตารางที่ 5.3 Traffic features

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

ตารางที่ 5.3 (ต่อ)

Feature name	Description
error_rate	RST error rate
srv_error_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 ประกอบด้วย 10 คุณลักษณะดังนี้

ตารางที่ 5.4 Host based features

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

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับ... อนุญาตให้นำไปใช้ประโยชน์ด้านการศึกษา

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

5.2.2 การเตรียมข้อมูลอินพุต

ในการดำเนินงานวิจัยนี้ผู้ดำเนินงานวิจัยได้เตรียมข้อมูลอินพุตสำหรับระบบ เพื่อใช้ในการทดลอง โดยข้อมูลที่ใช้ในการทดลองกำหนดบนพื้นฐานของประเภทการบุกรุกแบบ Denial of service กับ Probing และพฤติกรรมปกติรวม 11 ประเภท เนื่องจากเป็นกลุ่มข้อมูลที่มีจำนวนมากที่สุด มีทั้งสิ้น 476,945 เรคคอร์ด ดังแสดงในตารางที่ 5.5 และการจำแนกพฤติกรรมทั้งสามจะใช้เพียง 28 คุณลักษณะ [10] ประกอบด้วย 3 กลุ่ม คือ Basic features 9 คุณลักษณะ Time-based features 9 คุณลักษณะและ Host-based features 10 คุณลักษณะ

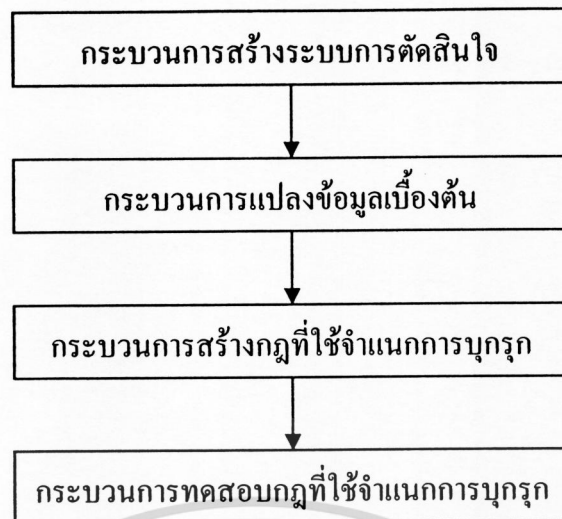
ตารางที่ 5.5 แสดงกลุ่มข้อมูลที่ใช้ทดลอง

Attack Type	จำนวนข้อมูล (เรคคอร์ด)
normal	81,548
neptune	107,201
smurf	280,790
back	1,925
satan	1,589
ipsweep	1,341
portsweep	1,057
teardrop	979
pod	264
land	20
nmap	231
รวม	476,945

5.3 การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่

การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่แสดงดังรูปที่ 5.2 ซึ่งประกอบด้วย 4 ขั้นตอน คือขั้นตอนที่หนึ่งกระบวนการสร้างระบบการตัดสินใจเป็นขั้นตอนที่นำข้อมูลที่ใช้ในการสอนและทดสอบระบบมาสร้างระบบการตัดสินใจ (decision system) ขั้นตอนที่สองกระบวนการแปลงข้อมูลเบื้องต้นเป็นขั้นตอนการแบ่งช่วงลักษณะเด่นที่มีความต่อเนื่องโดยใช้ฟิชชีเซต ขั้นตอนที่สามกระบวนการสร้างกฎที่ใช้จำแนกการบุกรุก (intrusion rules) เป็นขั้นตอนการสร้างกฎที่ใช้จำแนกการบุกรุกโดยใช้ราฟฟิซซี่ และขั้นตอนที่สี่กระบวนการทดสอบกฎที่ใช้จำแนกการบุกรุกเป็นขั้นตอนทดสอบประสิทธิภาพของกฎที่ใช้จำแนกการบุกรุกโดยการวัดค่า

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
Detection rate และ False negative rate คือ ไปจะกล่าวถึงรายละเอียดในแต่ละขั้นตอน
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 5.2 การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่

5.3.1 กระบวนการสร้างระบบการตัดสินใจ

ในการสร้างระบบการตัดสินใจ เพื่อใช้ในระบบตรวจจับการบุกรุกผ่านทางเครือข่ายโดยใช้ราฟฟิซซี่ จะดำเนินการโดยนำข้อมูลที่ใช้ในการสอนและทดสอบระบบมาสร้างระบบการตัดสินใจที่ถูกนิยามโดยสมการต่อไปนี้

$$S = (U, A \cup \{d\}) \quad (5.1)$$

โดยที่ U คือเซตของวัตถุ
 A คือเซตของคุณลักษณะที่ใช้ในการสร้างเงื่อนไข
 $d \notin A$ คือคุณลักษณะที่ใช้ในการตัดสินใจ และ $A \cap \{d\} = \emptyset$

ข้อมูลอินพุตที่ใช้ในงานวิจัยนี้ ประกอบด้วยคุณลักษณะที่ไม่มีความต่อเนื่องของข้อมูล ได้แก่ ชนิดของโปรโตคอล (protocol) บริการของเน็ตเวิร์คที่ทำงานบนเครื่องปลายทาง (service) และสถานะปกติหรือผิดพลาด (flag) เป็นต้น และคุณลักษณะที่มีความต่อเนื่องของข้อมูล ได้แก่ ระยะเวลาในการเชื่อมต่อ (duration) จำนวนข้อมูลที่ส่งจากเครื่องต้นทางไปยังเครื่องปลายทาง (src_bytes) จำนวนข้อมูลที่ส่งจากเครื่องปลายทางไปยังเครื่องต้นทาง (dst_bytes) จำนวนของการเชื่อมต่อที่มีโฮสต์เหมือนกันใน 2 นาทีที่ผ่านมา (count) จำนวนของการเชื่อมต่อที่มีบริการเหมือนกันใน 2 นาทีที่ผ่านมา (srv_count) เปอร์เซ็นต์ของการเชื่อมต่อที่มีบริการเหมือนกัน (same_srv_rate) และ เปอร์เซ็นต์ของการเชื่อมต่อที่มีบริการไม่เหมือนกัน (diff_srv_rate) เป็นต้น โดยตัวอย่างข้อมูลแสดงดังตารางที่ใช้ 5.6. ซึ่งประกอบด้วยคุณลักษณะที่มีความต่อเนื่อง และการคำนวณคุณลักษณะที่ไม่มีความต่อเนื่องรวมทั้งหมด 28 คุณลักษณะ อิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 5.6 แสดงตัวอย่างข้อมูล

คุณลักษณะ	ตัวอย่างข้อมูล				
	0	0	0	1	0
duration	0	0	0	1	0
protocol_type	tcp	tcp	tcp	tcp	icmp
service	http	private	http	private	ecr_i
flag	SF	S0	SF	RSTR	SF
src_bytes	54540	0	239	0	1032
dst_bytes	8314	0	486	0	0
land	0	0	0	0	0
wrong_fragment	0	0	0	0	0
urgent	0	0	0	0	0
count	3	16	8	18	511
srv_count	3	15	8	2	511
serror_rate	0.00	0.94	0.00	0.00	0.00
srv_serror_rate	0.00	1.00	0.00	0.00	0.00
error_rate	0.00	0.00	0.00	1.00	0.00
srv_error_rate	0.00	0.00	0.00	1.00	0.00
same_srv_rate	1.00	0.94	1.00	0.11	1.00
diff_srv_rate	0.00	0.12	0.00	0.50	0.00
srv_diff_host_rate	0.00	0.00	0.00	0.00	0.00
dst_host_count	255	15	19	188	158
dst_host_srv_count	255	16	19	2	13
dst_host_same_srv_rate	1.00	1.00	1.00	0.01	0.08
dst_host_diff_srv_rate	0.00	0.00	0.00	0.06	0.02
dst_host_same_src_port_rate	0.00	0.07	0.05	0.10	0.08
dst_host_srv_diff_host_rate	0.00	0.12	0.00	0.00	0.00
dst_host_serror_rate	0.00	1.00	0.00	0.01	0.00
dst_host_srv_serror_rate	0.00	0.94	0.00	0.00	0.00
dst_host_error_rate	0.05	0.00	0.00	0.36	0.00
dst_host_srv_error_rate	0.05	0.00	0.00	1.00	0.00
attack	back	neptune	normal	portsweep	smurf

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 5.6 จะแสดงตัวอย่างของข้อมูลที่ใช้ในงานวิจัยนี้ ซึ่งในตารางจะแสดงตัวอย่างข้อมูลของพฤติกรรมการบุกรุกประเภท back, neptune, portsweep, smurf, และพฤติกรรมปกติ normal โดย back เป็นพฤติกรรมการบุกรุกที่ทำให้เซิร์ฟเวอร์ทำงานช้า neptune เป็นพฤติกรรมการบุกรุกที่ทำให้เครื่องเป้าหมายไม่สามารถให้บริการได้ portsweep เป็นพฤติกรรมการบุกรุกที่สแกนหาโฮสต์ที่อยู่ในเครือข่ายเพื่อช่องโหว่และ smurf เป็นพฤติกรรมการบุกรุกที่ทำให้ระบบเครือข่ายหยุดทำงาน

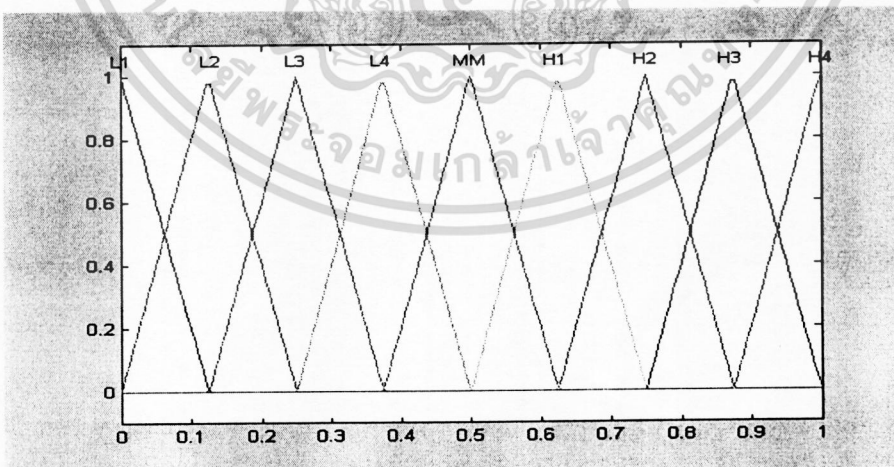
5.3.2 กระบวนการแปลงข้อมูลเบื้องต้น

ในกระบวนการแปลงข้อมูลเบื้องต้น จะทำการแปลงชนิดข้อมูลของคุณลักษณะที่มีลักษณะต่อเนื่อง (continuous) ให้เป็นช่วง (discrete) [12] โดยอันดับแรกจะทำการแปลงค่าของแต่ละคุณลักษณะให้มีค่าอยู่ระหว่าง 0 ถึง 1 โดยใช้สมการต่อไปนี้

$$x = \frac{x - MIN}{MAX - MIN} \quad (5.2)$$

โดยที่ x คือค่าของคุณลักษณะ
 MIN คือค่าต่ำที่สุดของคุณลักษณะ
 MAX คือค่าสูงที่สุดของคุณลักษณะ

จากนั้นทำการแปลงค่าคุณลักษณะที่มีลักษณะต่อเนื่องให้เป็นช่วงที่มีระยะห่างเท่าๆกัน โดยใช้ 9 fuzzy spaces ดังแสดงในรูปที่ 5.3



รูปที่ 5.3 Fuzzy spaces ของแต่ละคุณลักษณะ

และจากรูปที่ 5.3 ฟังก์ชันความเป็นสมาชิกที่เลือกจะเป็นแบบรูปสามเหลี่ยม เพราะเป็นฟังก์ชันความเป็นสมาชิกที่เหมาะสมสำหรับการแบ่งข้อมูลออกเป็นช่วงๆ และใช้ในการจำแนกวัตถุที่ดีที่สุด ซึ่งไม่สามารถคำนวณได้โดยสมการต่อไปนี้ เนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$$\text{triangle}(x; a, b, c) = \begin{cases} 0 & x \leq a \\ x - a & a \leq x \leq b \\ b - a & b \leq x \leq c \\ c - x & c \leq x \\ 0 & \end{cases} \quad (5.3)$$

หลังดำเนินการแบ่งช่วงคุณลักษณะจะได้ข้อมูลที่มีลักษณะดังที่แสดงในตารางที่ 5.7

ตารางที่ 5.7 แสดงตัวอย่างข้อมูลหลังผ่านการแปลงข้อมูลเบื้องต้น

คุณลักษณะ	ค่าข้อมูลก่อนการแปลงข้อมูลเบื้องต้น	ค่าข้อมูลหลังการแปลงข้อมูลเบื้องต้น
duration	0	L1
protocol_type	tcp	tcp
service	http	http
flag	SF	SF
src_bytes	239	L1L2
dst_bytes	486	L1L2
land	0	L1
wrong_fragment	0	L1
urgent	0	L1
count	8	L1L2
srv_count	8	L1L2
serror_rate	0.00	L1
srv_serror_rate	0.00	L1
rerror_rate	0.00	L1
srv_rerror_rate	0.00	L1
same_srv_rate	1.00	H4
diff_srv_rate	0.00	L1
srv_diff_host_rate	0.00	L1
dst_host_count	19	L2L1
dst_host_srv_count	19	L2L1
dst_host_same_srv_rate	1.00	H4
dst_host_diff_srv_rate	0.00	L1

ตารางที่ 5.7 (ต่อ)

คุณลักษณะ	ค่าข้อมูลก่อนการ แปลงข้อมูลเบื้องต้น	ค่าข้อมูลหลังการ แปลงข้อมูลเบื้องต้น
dst_host_same_src_port_rate	0.05	L1L2
dst_host_srv_diff_host_rate	0.00	L1
dst_host_serror_rate	0.00	L1
dst_host_srv_serror_rate	0.00	L1
dst_host_rerror_rate	0.00	L1
dst_host_srv_rerror_rate	0.00	L1
attack	normal	normal

ตารางที่ 5.7 แสดงการดำเนินการแปลงข้อมูลเบื้องต้นกับข้อมูลที่เป็นของพฤติกรรมปกติ ซึ่งจะทำให้แปลงชนิดข้อมูลของคุณลักษณะที่มีลักษณะต่อเนื่องให้เป็นข้อมูลที่มีลักษณะเป็นช่วง เช่น คุณลักษณะ Duration ค่าก่อนทำการแปลงเท่ากับ 0 วินาทีหลังทำการแปลงจะมีค่าเท่ากับ L1 Count ค่าก่อนทำการแปลงเท่ากับ 8 คอนเน็คชันหลังทำการแปลงจะมีค่าเท่ากับ L1L2 และ Same_srv_rate ค่าก่อนทำการแปลงเท่ากับ 100 เปอร์เซ็นต์หลังทำการแปลงจะมีเท่ากับ H4 เป็นต้น

5.3.3 กระบวนการสร้างกฎที่ใช้จำแนกการบุกรุก

ในกระบวนการนี้ จะทำการสร้างกฎที่ใช้จำแนกการบุกรุก (intrusion rules) โดยใช้ reducts หรือซับเซตที่มีคุณลักษณะน้อยที่สุดที่ยังคงสามารถจำแนกวัตถุได้ที่คำนวณหาได้จาก อัลกอริทึมลดคุณลักษณะ โดยในงานวิจัยที่ได้นำเสนอจะใช้อัลกอริทึม discernibility matrix และ Boolean function ที่มีขั้นตอนดังนี้

1. หา discernibility matrix M โดยสมาชิก m_{ij} สามารถหาได้จากสมการต่อไปนี้

$$(m_{ij}) = \begin{cases} \{a \in A : a(x_i) \neq a(x_j)\}, & D(x_i) \neq D(x_j) \\ \phi, & D(x_i) = D(x_j) \end{cases} \quad (5.4)$$

2. หา discernibility function L_{ij} จาก discernibility matrix M โดยใช้สมการต่อไปนี้

$$L_{ij} = \bigvee_{a_i \in m_{ij}} a_i, \quad m_{ij} \neq \phi \quad (5.5)$$

3. หา discernibility function L โดยใช้สมการต่อไปนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่นิยมนำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

$$L = \bigwedge_{m_j \neq \phi} L_{ij} \quad (5.6)$$

4. ทำการลดรูป discernibility function L ที่ได้จากข้อ 3 โดยใช้ absorption laws ต่อไปนี้

$$\begin{aligned} a \vee (a \wedge b) &= a \\ a \wedge (a \vee b) &= a \end{aligned} \quad (5.7)$$

5. ทำการแปลง discernibility function ที่ได้จากข้อ 4 ซึ่งอยู่ในรูปแบบ DNF (Disjunctive Normal Form) ให้อยู่ในรูปแบบ CNF (Conjunctive Normal Form) โดยใช้ distribution law ต่อไปนี้

$$\begin{aligned} a \wedge (a \vee b) &= (a \wedge b) \vee (a \wedge c) \\ a \vee (a \wedge b) &= (a \vee b) \wedge (a \vee c) \end{aligned} \quad (5.8)$$

6. ทำการลดรูป discernibility function ที่ได้จากข้อ 5 โดยใช้ absorption laws ที่นิยามโดยสมการที่ 5.7 จากนั้นก็จะได้ reducts หรือซัพเซตที่มีคุณลักษณะน้อยที่สุดที่ยังคงสามารถจำแนกวัตถุได้
7. กำหนดหาค่าความขึ้นต่อกันของคุณลักษณะที่อยู่ในแต่ละ reducts โดยใช้สมการต่อไปนี้

$$k = \frac{|POS_C(D)|}{|U|} \quad (5.9)$$

8. เลือก reduct หรือซัพเซตของคุณลักษณะที่มีค่าความขึ้นต่อกันสูงที่สุดมาสร้างกฎที่ใช้จำแนกการบุกรุก ซึ่งมีรูปแบบดังนี้สมการต่อไปนี้

$$(a_{i_1} = v_1) \wedge \dots \wedge (a_{i_m} = v_m) \Rightarrow (d = v) \quad (5.10)$$

เมื่อดำเนินการตามอัลกอริทึมเสร็จ จะได้ reducts ซึ่งเป็นซัพเซตของคุณลักษณะที่จำเป็นต่อการจำแนกพฤติกรรมบุกรุกและพฤติกรรมปกติดังตารางที่ 5.8 และเมื่อนำ reduct ที่มีค่าการขึ้นต่อกันของคุณลักษณะสูงที่สุดและมีคุณลักษณะน้อยที่สุดมาสร้างกฎที่ใช้จำแนกการบุกรุกจะได้กฎที่ใช้จำแนกการบุกรุกดังตารางที่ 5.9

ตารางที่ 5.8 แสดงตัวอย่าง reducts

จำนวน คุณลักษณะ	ค่าการขึ้น ต่อกันของ คุณลักษณะ	Reducts
14	0.998	duration, wrong, count, srv_error, srv_error, srv_diff_host, dst_host_count, dst_host_srv_count, dst_host_diff_srv, dst_host_srv_diff, dst_host_error ,dst_host_srv_error, dst_host_error, service
15	0.998	duration, wrong, count, srv_error, srv_error, srv_diff_host, dst_host_count, dst_host_srv_count, dst_host_diff_srv, dst_host_same_src_port, dst_host_srv_diff, dst_host_error, dst_host_srv_error, dst_host_srv_error, service
...

ตารางที่ 5.9 แสดงตัวอย่างกฎที่ใช้จำแนกการบุกรุก

Condition	Action
(duration=L1)&(wrong=L1)&(count=L1L2)&(srv_error=H3H2)& (srv_error=L1)&(srv_diff_host=L1)&(dst_host_count=L1L2)& (dst_host_srv_count=L1L2)&(dst_host_diff_srv=L1)&(dst_host_srv_diff=L4L3)& (dst_host_error=H4)&(dst_host_srv_error=H3H2)&(dst_host_error=L1)& (service=private)	neptune
(duration=L1)&(wrong=L1)&(count=L1L2)&(srv_error=L1)&(srv_error=L1)& (srv_diff_host=L1)&(dst_host_count=L1L2)&(dst_host_srv_count=L1L2)& (dst_host_diff_srv=L1)&(dst_host_srv_diff=L1)&(dst_host_error=L1)& (dst_host_srv_error=L1)&(dst_host_error=L1))&(service=http)	normal
(duration=L1L2)&(wrong=L1)&(count=L1L2)&(srv_error=L1)& (srv_error=H4)&(srv_diff_host=L1)&(dst_host_count=H2H1)& (dst_host_srv_count=L1L2)&(dst_host_diff_srv=L1L2)&(dst_host_srv_diff=L1) &(dst_host_error=L1L2)&(dst_host_srv_error=L1)&(dst_host_error=L4L3)& (service=private)	portsweep
...	...

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

5.3.4 กระบวนการทดสอบกฎที่ใช้จำแนกการบุกรุก

ในกระบวนการนี้ จะทำการทดสอบกฎที่ใช้จำแนกการบุกรุกกับกลุ่มข้อมูลที่ใช้ทดสอบ ซึ่งจะเน้นเฉพาะค่า Detection rate และ False negative rate เท่านั้น

Detection rate คือค่าเปอร์เซ็นต์ของจำนวนการบุกรุกที่ตรวจจับได้ถูกต้องกับจำนวนของการบุกรุกที่ตรวจจับได้ทั้งหมดที่ถูกระบุประเภทการบุกรุกโดยกฎที่ใช้จำแนกการบุกรุก คำนวณได้ดังนี้

$$\%Detected = \left(\frac{N_{corrected}}{N_{corrected} + N_{missed}} \right) \times 100 \quad (5.11)$$

โดย $\%Detected$ คือค่า Detection rate

$N_{corrected}$ คือจำนวนข้อมูลทดสอบที่กฎที่ใช้จำแนกการบุกรุกจำแนกได้ถูกต้อง
 N_{missed} คือจำนวนข้อมูลทดสอบที่กฎที่ใช้จำแนกการบุกรุกจำแนกผิดพลาด

False negative rate คือค่าเปอร์เซ็นต์ของจำนวนการบุกรุกที่กฎที่ใช้จำแนกพฤติกรรมผิดปกติจำแนกผิดพลาดกับจำนวนพฤติกรรมปกติที่จำแนกได้ถูกต้อง คำนวณได้ดังนี้

$$\%FalseNegative = \left(\frac{N_{false}}{N_{normal}} \right) \times 100 \quad (5.12)$$

โดย $\%FalseNegative$ คือค่า False negative rate

N_{false} คือจำนวนของการบุกรุกที่กฎที่ใช้จำแนกพฤติกรรมผิดปกติจำแนกผิดพลาด
 N_{normal} คือจำนวนพฤติกรรมปกติที่กฎที่ใช้จำแนกพฤติกรรมปกติจำแนกได้ถูกต้อง

5.4 สรุป

ในบทนี้ผู้วิจัยได้นำเสนอวิธีการตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่ และใช้ข้อมูลที่ได้จาก KDD Cup 1999 [9] ในการทดลอง โดยทดลองกับข้อมูลการบุกรุกแบบ Denial of service (DoS) Probing และพฤติกรรมปกติรวม 11 ประเภท อ้างอิงจาก KDD Cup 1999 [9]

การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่ ประกอบด้วย 4 ขั้นตอน คือ ขั้นตอนหนึ่งเป็นการสร้างระบบการตัดสินใจที่นำข้อมูลที่ใช้ในการสอนและทดสอบระบบมาสร้างระบบการตัดสินใจ ขั้นตอนที่สองเป็นการแบ่งช่วงลักษณะเด่นที่มีความต่อเนื่องโดยใช้ฟัซซี่เซต ขั้นตอนที่สามเป็นการสร้างกฎที่ใช้จำแนกการบุกรุกโดยใช้ราฟฟิซซี่ และขั้นตอนที่สี่ทดสอบประสิทธิภาพของกฎที่ใช้จำแนกการบุกรุกโดยการวัดค่า Detection rate และ False negative rate

เพื่อพิสูจน์ว่าวิธีการตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟิซที่สามารถนำไปประยุกต์ใช้ได้
อย่างมีประสิทธิภาพ



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 6

ผลการทดลอง

ในบทนี้จะกล่าวถึงการนำระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิชซึ่งจากบทที่ 5 มาทดลอง และทำการเปรียบเทียบประสิทธิภาพในการตรวจจับการบุกรุกกับระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิวรอลเน็ตเวิร์คและดีซิชั่นทรีตามลำดับ โดยในการเปรียบเทียบจะใช้การวัดค่า Detection rate และค่า False negative rate ตามลำดับ

6.1 การเตรียมเครื่องมือและข้อมูลที่ใช้ในการทดลอง

6.1.1 เครื่องมือที่ใช้ในการทดลอง

ในงานวิจัยนี้ใช้เครื่องประมวลผลเป็นเครื่องไมโครคอมพิวเตอร์ที่ใช้ซีพียูเพนเทียมของบริษัทอินเทลรุ่น P4 2.4 GHz หน่วยความจำ 768 MB และโปรแกรมที่ใช้ในการเขียนโปรแกรมคือโปรแกรม MATLAB เวอร์ชัน 6.5 ซึ่งได้ใช้ในส่วนที่เกี่ยวกับการเขียนโปรแกรม นิวรอลเน็ตเวิร์คทูลบ็อกซ์ (Neural Network Toolbox) และดีซิชั่นทรีทูลบ็อกซ์ (Decision Tree Toolbox)

6.1.2 การเตรียมข้อมูลที่ใช้ในการทดลอง

ข้อมูลที่นำมาทดลองในงานวิจัยนี้เป็นชุดข้อมูลของ KDD Cup 1999 โดยการทดลองได้กำหนดบนพื้นฐานของประเภทการบุกรุกแบบ Denial of service (DoS) Probing และพฤติกรรมปกติรวม 11 ประเภท เนื่องจากเป็นกลุ่มข้อมูลที่มีจำนวนมากที่สุด มีทั้งสิ้น 476,945 เรคคอร์ด ดังแสดงในตารางที่ 6.1 และการจำแนกพฤติกรรมทั้งสามจะใช้เพียง 28 คุณลักษณะ [10] ประกอบด้วย 3 กลุ่ม คือ Basic features 9 คุณลักษณะ Time-based features 9 คุณลักษณะและ Host-based features 10 คุณลักษณะ ดังแสดงในตารางที่ 6.2 เพื่อให้ข้อมูลมีความหลากหลายและครอบคลุมในงานวิจัยนี้จึงได้สร้างข้อมูลที่ใช้ในการสอนและทดสอบระบบขึ้นจำนวน 5 ชุด ซึ่งในละชุดจะแบ่งข้อมูลที่ใช้ในการสอนและทดสอบระบบด้วยการสุ่มซ้ำๆ กันชุดละ 3 ครั้ง โดยข้อมูลชุดที่ 1 จะสุ่มข้อมูลที่ใช้สอนระบบ 20% และทดสอบระบบ 80 % ข้อมูลชุดที่ 2 จะสุ่มข้อมูลที่ใช้สอนระบบ 40% และทดสอบระบบ 60 % ข้อมูลชุดที่ 3 จะสุ่มข้อมูลที่ใช้สอนระบบ 50% และทดสอบระบบ 50 % ข้อมูลชุดที่ 4 จะสุ่มข้อมูลที่ใช้สอนระบบ 60% และทดสอบระบบ 40 % และข้อมูลชุดที่ 5 จะสุ่มข้อมูลที่ใช้สอนระบบ 80% และทดสอบระบบ 20 % ตามลำดับ ดังตารางที่ 6.3

ตารางที่ 6.1 แสดงกลุ่มข้อมูลที่ใช้ทดลอง

Attack Type	จำนวนข้อมูล (เรคคอร์ด)
normal	81,548
neptune	107,201
smurf	280,790
back	1,925
satan	1,589
ipsweep	1,341
portsweep	1,057
teardrop	979
pod	264
land	20
nmap	231
รวม	476,945

ตารางที่ 6.2 แสดงตัวอย่างข้อมูลที่ใช้ในการสอนและทดสอบระบบ

คุณลักษณะ	ตัวอย่างข้อมูล				
duration	0	0	0	1	0
protocol_type	tcp	tcp	tcp	tcp	icmp
service	http	private	http	private	ecr_i
flag	SF	S0	SF	RSTR	SF
src_bytes	54540	0	239	0	1032
dst_bytes	8314	0	486	0	0
Land	0	0	0	0	0
wrong_fragment	0	0	0	0	0
urgent	0	0	0	0	0
count	3	16	8	18	511
srv_count	3	15	8	2	511
serror_rate	0.00	0.94	0.00	0.00	0.00
srv_serror_rate	0.00	1.00	0.00	0.00	0.00
error_rate	0.00	0.00	0.00	1.00	0.00
srv_rerror_rate	0.00	0.00	0.00	1.00	0.00

ตารางที่ 6.2 (ต่อ)

คุณลักษณะ	ตัวอย่างข้อมูล				
	1.00	0.94	1.00	0.11	1.00
same_srv_rate	1.00	0.94	1.00	0.11	1.00
diff_srv_rate	0.00	0.12	0.00	0.50	0.00
srv_diff_host_rate	0.00	0.00	0.00	0.00	0.00
dst_host_count	255	15	19	188	158
dst_host_srv_count	255	16	19	2	13
dst_host_same_srv_rate	1.00	1.00	1.00	0.01	0.08
dst_host_diff_srv_rate	0.00	0.00	0.00	0.06	0.02
dst_host_same_src_port_rate	0.00	0.07	0.05	0.10	0.08
dst_host_srv_diff_host_rate	0.00	0.12	0.00	0.00	0.00
dst_host_serror_rate	0.00	1.00	0.00	0.01	0.00
dst_host_srv_serror_rate	0.00	0.94	0.00	0.00	0.00
dst_host_error_rate	0.05	0.00	0.00	0.36	0.00
dst_host_srv_error_rate	0.05	0.00	0.00	1.00	0.00
attack	back	neptune	normal	portsweep	smurf

ตารางที่ 6.3 แสดงรูปแบบการแบ่งข้อมูลที่ใช้ทดลอง

รูปแบบ	ข้อมูลที่ให้สอน	ข้อมูลที่ให้ทดสอบ	รวม
	(เรคคอร์ด)	(เรคคอร์ด)	(เรคคอร์ด)
20/80	95,394	381,551	476,945
40/60	190,783	286,192	476,945
50/50	238,475	238,470	476,945
60/40	286,165	190,780	476,945
80/20	381,553	95,392	476,945

ในงานวิจัยนี้ได้นำวิธีตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิรอลเน็ตเวิร์คและดีซีชันทรีมาทดลองเปรียบเทียบประสิทธิภาพกับวิธีตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่ที่นำเสนอ ซึ่งวิธีทั้งสองจะประมวลผลข้อมูลที่มีลักษณะเป็นตัวเลข ดังนั้นต้องทำการแปลงคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์ให้เป็นตัวเลขก่อน ซึ่งคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์ประกอบด้วย Protocol feature Service feature และ Flag feature โดยมีรายละเอียดดังต่อไปนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 6.4 แสดงการกำหนดค่าตัวเลขแทน Protocol feature

Value	Assigned
tcp	0
udp	1
icmp	2

ตารางที่ 6.5 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
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
nntp	31	bgp	55
http_443	32	Z39_50	56
exec	33	tim_i	64

ตารางที่ 6.5 (ต่อ)

Value	Assigned	Value	Assigned
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
systat	42	tim_i	64
supdup	43	red_j	65

ตารางที่ 6.6 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		

6.2 การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิชซี่

การทดลองระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิชซี่ในงานวิจัยนี้ จะใช้ชุดข้อมูลของ KDD Cup 1999 [9] สำหรับสอนและทดสอบระบบ (ข้อมูลในข้อ 6.1.2) ซึ่งได้มีการเตรียมข้อมูลสำหรับใช้สอนและทดสอบระบบจำนวน 5 ชุดดังตารางที่ 6.3 การทดลองจะทำการทดสอบระบบกับข้อมูลทั้ง 5 ชุดและนำค่า Detection rate และ False negative rate ของข้อมูลแต่ละชุดมาคำนวณหาค่าเฉลี่ย โดยมีขั้นตอนการทดลองดังนี้

- นำคุณลักษณะที่มีความต่อเนื่องของข้อมูลที่ใช้สอนและทดสอบระบบแต่ละชุดไปแปลงให้มีค่าอยู่ระหว่าง 0 ถึง 1 โดยใช้สมการ 5.2 ในบทที่ 5 หัวข้อ 5.3.2
- นำคุณลักษณะที่มีความต่อเนื่องจากข้อ 1. ไปแปลงให้เป็นช่วงโดยใช้ fuzzy spaces ที่แสดงในรูปที่ 5.3 ในบทที่ 5 หัวข้อ 5.3.2
- นำข้อมูลที่ใช้สอนระบบแต่ละชุดไปสอนระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิชซี่ที่ใช้อัลกอริทึมในบทที่ 5 หัวข้อ 5.3.3

4. นำข้อมูลที่ใช้ทดสอบระบบแต่ละชุดไปทดสอบระบบ

ต่อไปจะอธิบายรายละเอียดในแต่ละขั้นตอน

1. นำคุณลักษณะที่มีความต่อเนื่องของข้อมูลที่ใช้สอนและทดสอบระบบแต่ละชุดไปแปลงให้มีค่าอยู่ระหว่าง 0 ถึง 1 โดยใช้สมการ 5.2 ในบทที่ 5 หัวข้อ 5.3.2 ซึ่งหลังจากทำการแปลงค่าจะได้ผลดังตารางที่ 6.7

ตารางที่ 6.7 แสดงตัวอย่างข้อมูลที่ถูกแปลงค่าคุณลักษณะที่มีความต่อเนื่องให้มีค่าอยู่ระหว่าง 0

ถึง 1

คุณลักษณะ	ตัวอย่างข้อมูล				
duration	0.00	0.00	0.00	0.00	0.00
protocol_type	tcp	Tcp	tcp	tcp	icmp
service	http	Private	http	private	ecr_i
flag	SF	S0	SF	RSTR	SF
src_bytes	0.00	0.00	0.00	0.00	0.00
dst_bytes	0.00	0.00	0.00	0.00	0.00
Land	0.00	0.00	0.00	0.00	0.00
wrong_fragment	0.00	0.00	0.00	0.00	0.00
urgent	0.00	0.00	0.00	0.00	0.00
count	0.00	0.03	0.01	0.03	1.00
srv_count	0.00	0.03	0.01	0.00	1.00
serror_rate	0.00	0.94	0.00	0.00	0.00
srv_serror_rate	0.00	1.00	0.00	0.00	0.00
rerror_rate	0.00	0.00	0.00	1.00	0.00
srv_rerror_rate	0.00	0.00	0.00	1.00	0.00
same_srv_rate	1.00	0.94	1.00	0.11	1.00
diff_srv_rate	0.00	0.12	0.00	0.50	0.00
srv_diff_host_rate	0.00	0.00	0.00	0.00	0.00
dst_host_count	1.00	0.06	0.07	0.74	0.62
dst_host_srv_count	1.00	0.06	0.07	0.00	0.05
dst_host_same_srv_rate	1.00	1.00	1.00	0.01	0.08
dst_host_diff_srv_rate	0.00	0.00	0.00	0.06	0.02
dst_host_same_src_port_rate	0.00	0.07	0.05	0.10	0.08

ตารางที่ 6.7 (ต่อ)

คุณลักษณะ	ตัวอย่างข้อมูล				
	0.00	0.12	0.00	0.00	0.00
dst_host_srv_diff_host_rate	0.00	0.12	0.00	0.00	0.00
dst_host_serror_rate	0.00	1.00	0.00	0.01	0.00
dst_host_srv_serror_rate	0.00	0.94	0.00	0.00	0.00
dst_host_rerror_rate	0.05	0.00	0.00	0.36	0.00
dst_host_srv_rerror_rate	0.05	0.00	0.00	1.00	0.00
attack	back	neptune	normal	portsweep	smurf

ตารางที่ 6.7 แสดงผลการแปลงค่าคุณลักษณะที่มีความต่อเนื่องของข้อมูลที่ใช้สอนและทดสอบระบบแต่ละชุดให้มีค่าอยู่ระหว่าง 0 ถึง 1 ซึ่งคุณลักษณะที่มีถูกแปลงค่าประกอบด้วยคุณลักษณะ duration count และ same_srv_rate เป็นต้น

2. นำคุณลักษณะที่มีความต่อเนื่องจากข้อ 1. ไปแปลงให้เป็นช่วงโดยใช้ fuzzy spaces ที่แสดงในรูปที่ 5.3 ในบทที่ 5 หัวข้อ 5.3.2 ซึ่งหลังจากทำการแปลงค่าจะได้ผลดังตารางที่ 6.8

ตารางที่ 6.8 แสดงตัวอย่างข้อมูลที่ถูกแปลงค่าคุณลักษณะที่มีความต่อเนื่องให้เป็นช่วงโดยใช้ fuzzy spaces

คุณลักษณะ	ตัวอย่างข้อมูล				
	L1	L1	L1	L1L2	L1
duration	L1	L1	L1	L1L2	L1
protocol_type	tcp	tcp	tcp	tcp	icmp
service	http	private	http	private	ecr_i
flag	SF	S0	SF	RSTR	SF
src_bytes	L1L2	L1	L1L2	L1	L1L2
dst_bytes	L1L2	L1	L1L2	L1	L1
Land	L1	L1	L1	L1	L1
wrong_fragment	L1	L1	L1	L1	L1
urgent	L1	L1	L1	L1	L1
count	L1L2	L1L2	L1L2	L1L2	H4
srv_count	L1L2	L1L2	L1L2	L1L2	H4
serror_rate	L1	H4H3	L1	L1	L1
srv_serror_rate	L1	H4	L1	L1	L1
rerror_rate	L1	L1	L1	H4	L1

ตารางที่ 6.9 (ต่อ)

ชุดข้อมูล (สอน/ทดสอบ)	จำนวน คุณลักษณะ	ค่าการขึ้นต่อกัน ของคุณลักษณะ	จำนวนกฎที่ใช้จำแนก พฤติกรรมการบุกรุก
50/50	18	0.997	13,144
60/40	18	0.997	14,568
80/20	16	0.997	16,342

ตารางที่ 6.9 แสดงผลการสอนระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้กราฟฟิซซี่ด้วยข้อมูลที่ใช้สอนระบบแต่ละชุดข้อมูล จากตารางจำนวนคุณลักษณะและค่าการขึ้นต่อกันของคุณลักษณะไม่ได้ขึ้นกับจำนวนข้อมูลที่ใช้สอนแต่จะขึ้นกับความแตกต่างของข้อมูล และจำนวนของกฎที่ใช้จำแนกพฤติกรรมการบุกรุกจะเพิ่มขึ้นตามจำนวนข้อมูลที่ใช้สอนระบบ

ตารางที่ 6.10 แสดงตัวอย่างกฎที่ใช้จำแนกการบุกรุกที่ได้จากการเรียนรู้ของระบบ

Condition	Action
(duration=L1)&(wrong=L1)&(count=L1L2)&(srv_serror=H3H2)&(srv_rerror=L1)&(srv_diff_host=L1)&(dst_host_count=L1L2)&(dst_host_srv_count=L1L2)&(dst_host_diff_srv=L1)&(dst_host_srv_diff=L4L3)&(dst_host_serror=H4)&(dst_host_srv_serror=H3H2)&(dst_host_rerror=L1)&(service=private)	neptune
(duration=L1)&(wrong=L1)&(count=L1L2)&(srv_serror=L1)&(srv_rerror=L1)&(srv_diff_host=L1)&(dst_host_count=L1L2)&(dst_host_srv_count=L1L2)&(dst_host_diff_srv=L1)&(dst_host_srv_diff=L1)&(dst_host_serror=L1)&(dst_host_srv_serror=L1)&(dst_host_rerror=L1)&(service=http)	normal
(duration=L1L2)&(wrong=L1)&(count=L1L2)&(srv_serror=L1)&(srv_rerror=H4)&(srv_diff_host=L1)&(dst_host_count=H2H1)&(dst_host_srv_count=L1L2)&(dst_host_diff_srv=L1L2)&(dst_host_srv_diff=L1)&(dst_host_serror=L1L2)&(dst_host_srv_serror=L1)&(dst_host_rerror=L4L3)&(service=private)	portsweep
...	...

4. นำข้อมูลที่ใช้ทดสอบระบบแต่ละชุดไปทดสอบระบบ โดยผลการทดลองแสดงในตารางที่ 6.11 และ 6.12 ตามลำดับ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 6.11 Detection rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่

Attack Type	Detection Rate (%)					Average
	20/80	40/60	50/50	60/40	80/20	
back	98.04	97.63	96.04	98.15	97.20	97.41
ipsweep	99.67	99.05	99.83	99.19	98.80	99.31
land	100.00	100.00	100.00	100.00	100.00	100.00
neptune	100.00	100.00	100.00	100.00	100.00	100.00
nmap	95.46	97.78	98.28	97.73	96.15	97.08
normal	99.52	99.56	99.62	99.56	99.58	99.57
pod	100.00	100.00	99.06	100.00	100.00	99.81
portsweep	99.84	99.80	100.00	100.00	100.00	99.93
satan	99.90	100.00	100.00	100.00	100.00	99.98
smurf	100.00	100.00	100.00	99.99	100.00	100.00
teardrop	100.00	100.00	100.00	100.00	100.00	100.00
Average	99.31	99.44	99.35	99.51	99.24	99.37

ตารางที่ 6.11 แสดงผลการวัดค่า Detection rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่ จากตารางจะพบว่าระบบมีค่า Detection rate เฉลี่ยในการตรวจจับพฤติกรรมปกติและพฤติกรรมบุกรุกแต่ละประเภทของข้อมูลแต่ละชุดสูง โดยมีพฤติกรรมบุกรุกที่มีค่า Detection rate เท่ากับ 100.00% จำนวน 4 ประเภทได้แก่ land, neptune, smurf, และ teardrop เมื่อนำค่า Detection rate ของข้อมูลแต่ละชุดมาคำนวณหาค่าเฉลี่ยจะได้ค่า Detection rate เฉลี่ยเท่ากับ 99.37% ดังนั้นจึงสรุปได้ว่าระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่เป็นระบบที่มีเปอร์เซ็นต์ความแม่นยำในการตรวจจับพฤติกรรมบุกรุกสูง

ตารางที่ 6.12 False negative rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่

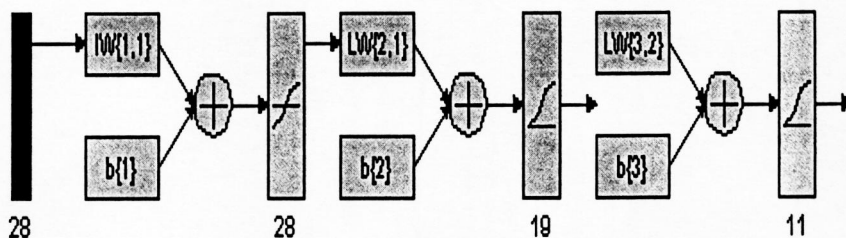
ชุดข้อมูล (สอน/ทดสอบ)	False negative rate (%)
20/80	0.062
40/60	0.071
50/50	0.113
60/40	0.056
80/20	0.094
Average	0.079

ตารางที่ 6.12 แสดงผลการวัดค่า False negative rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่ ซึ่งวัดจากการตรวจจับพฤติกรรมการบุกรุกที่ผิดพลาดของกฎที่ใช้จำแนกพฤติกรรมปกติ จากตารางจะพบว่าระบบมีค่า False negative rate ของข้อมูลแต่ละชุดต่ำกว่า 0.150% เมื่อนำค่า False negative rate ของข้อมูลแต่ละชุดมาคำนวณหาค่าเฉลี่ยจะได้ค่า False negative rate เฉลี่ยเท่ากับ 0.079% ดังนั้นจึงสรุปได้ว่าระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซี่เป็นระบบที่มีเปอร์เซ็นต์ความผิดพลาดในการตรวจจับพฤติกรรมการบุกรุกต่ำมาก

6.3 การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิวรอลเน็ตเวิร์ค

นิวรอลเน็ตเวิร์ค คือระบบที่มีการประมวลผลข้อมูลซึ่งรวมคุณสมบัติของไบโอลอจิกคอล นิวรอลเน็ตเวิร์คถูกพัฒนาขึ้นโดยโมเดลทางคณิตศาสตร์ของกระบวนการเรียนรู้ของมนุษย์ (เลียนแบบการทำงานของสมอง) และจะเรียนรู้จากชุดข้อมูลของชุดความรู้ทรงนิ่งเซต นิวรอลเน็ตเวิร์คประกอบด้วยหน่วยความจำจำนวนมากที่เรียกว่า นิวรอน (neurons) เซล (cells) หรือ โหนด (nodes) แต่ละนิวรอนจะต่อกันโดยคอนเน็คชันลิงค์ (connection link) ที่มีค่าน้ำหนักของมันอยู่ในแต่ละการเชื่อมต่อ โดยค่าน้ำหนักจะแสดงรายละเอียดที่เน็ตเวิร์คใช้ในการแก้ปัญหาโดยนิวรอลเน็ตเวิร์คถูกใช้ในการแก้ปัญหอย่างกว้างขวาง เช่น การเก็บและการเรียกข้อมูล การแยกประเภทของข้อมูล การเปลี่ยนจากรูปแบบของอินพุทให้อยู่ในรูปแบบของเอาต์พุท ความสามารถในการตรวจสอบรูปแบบของข้อมูลที่คล้ายคลึงกับความคิดของมนุษย์ เป็นต้น

การใช้นิวรอลเน็ตเวิร์คในงานวิจัยนี้มีจุดประสงค์เพื่อเปรียบเทียบประสิทธิภาพในการจำแนกพฤติกรรมการบุกรุกกับระบบตรวจจับการบุกรุกที่นำเสนอในงานวิจัยนี้ โดยนิวรอลเน็ตเวิร์คที่ใช้จะเป็นแบบ feed-forward backpropagation ซึ่งเป็นนิวรอลเน็ตเวิร์คแบบป้อนไปข้างหน้าชนิดหลายเลเยอร์ที่ใช้วิธีการเรียนรู้แบบแพร่ย้อนกลับ (Back propagation) ประกอบด้วย 3 เลเยอร์ คือ อินพุตเลเยอร์จำนวน 28 นิวรอน 1 เลเยอร์ ฮิดเดนเลเยอร์จำนวน 19 นิวรอน 1 เลเยอร์ เอาต์พุตเลเยอร์จำนวน 11 นิวรอน 1 เลเยอร์ ดังแสดงในรูปที่ 6.1 โดยในการทดลองจะใช้นิวรอลเน็ตเวิร์ค ทูลบ็อกซ์ที่อยู่ในโปรแกรม MATLAB เวอร์ชัน 6.5



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น **รูปที่ 6.1** โครงสร้าง BP network ของระบบตรวจจับการบุกรุก

สำหรับการทดลองจะใช้ข้อมูลชุดเดียวกันกับที่ใช้กับระบบที่นำเสนอในงานวิจัยนี้ (ข้อมูลในข้อ 6.1.2) ซึ่งได้มีการเตรียมข้อมูลสำหรับใช้สอนและทดสอบระบบ 5 ชุดดังตารางที่ 6.3 โดยจะทำการทดสอบระบบกับข้อมูลทั้ง 5 ชุดและนำค่า Detection rate และ False negative rate ของข้อมูลแต่ละชุดมาหาค่าเฉลี่ย โดยมีขั้นตอนการทดลองดังนี้

1. นำคุณลักษณะที่มีความต่อเนื่องของข้อมูลที่ใช้สอนและทดสอบระบบแต่ละชุดไปแปลงให้มีค่าอยู่ระหว่าง 0 ถึง 1 โดยใช้สมการ 5.2 ในบทที่ 5 หัวข้อ 5.3.2
2. นำคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์จากข้อ 1. ไปแปลงค่าเป็นตัวเลข
3. นำข้อมูลที่ใช้สอนระบบแต่ละชุดไปสอนระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิวรอลเน็ตเวิร์ค
4. นำข้อมูลที่ใช้ทดสอบระบบแต่ละชุดไปทดสอบระบบ

ต่อไปจะอธิบายรายละเอียดในแต่ละขั้นตอน

1. นำคุณลักษณะที่มีความต่อเนื่องของข้อมูลที่ใช้สอนและทดสอบระบบแต่ละชุดไปแปลงให้มีค่าอยู่ระหว่าง 0 ถึง 1 โดยใช้สมการ 5.2 ในบทที่ 5 หัวข้อ 5.3.2 ซึ่งหลังจากทำการแปลงค่าจะได้ผลดังตารางที่ 6.13

ตารางที่ 6.13 แสดงตัวอย่างข้อมูลที่ถูกแปลงค่าคุณลักษณะที่มีความต่อเนื่องให้มีค่าอยู่ระหว่าง 0 ถึง 1

คุณลักษณะ	ตัวอย่างข้อมูล				
	0.00	0.00	0.00	0.00	0.00
duration	0.00	0.00	0.00	0.00	0.00
protocol_type	tcp	tcp	tcp	tcp	icmp
service	http	private	http	private	ecr_i
flag	SF	S0	SF	RSTR	SF
src_bytes	0.00	0.00	0.00	0.00	0.00
dst_bytes	0.00	0.00	0.00	0.00	0.00
Land	0.00	0.00	0.00	0.00	0.00
wrong_fragment	0.00	0.00	0.00	0.00	0.00
urgent	0.00	0.00	0.00	0.00	0.00
count	0.00	0.03	0.01	0.03	1.00
srv_count	0.00	0.03	0.01	0.00	1.00
serror_rate	0.00	0.94	0.00	0.00	0.00
srv_serror_rate	0.00	1.00	0.00	0.00	0.00
rerror_rate	0.00	0.00	0.00	1.00	0.00
srv_rerror_rate	0.00	0.00	0.00	1.00	0.00

เอกสารที่สงวนไว้สำหรับการใช้งานในการศึกษานั้น ไม่อนุญาตให้นำไปเผยแพร่โดยไม่ได้รับอนุญาต
 ไม่ได้รับอนุญาตให้คัดลอกหรือทำซ้ำโดยไม่ได้รับอนุญาต และหากมีข้อสงสัยหรือต้องการข้อมูลเพิ่มเติม กรุณาติดต่อผู้จัดทำเอกสาร

ตารางที่ 6.13 (ต่อ)

คุณลักษณะ	ตัวอย่างข้อมูล				
	same_srv_rate	1.00	0.94	1.00	0.11
diff_srv_rate	0.00	0.12	0.00	0.50	0.00
srv_diff_host_rate	0.00	0.00	0.00	0.00	0.00
dst_host_count	1.00	0.06	0.07	0.74	0.62
dst_host_srv_count	1.00	0.06	0.07	0.00	0.05
dst_host_same_srv_rate	1.00	1.00	1.00	0.01	0.08
dst_host_diff_srv_rate	0.00	0.00	0.00	0.06	0.02
dst_host_same_src_port_rate	0.00	0.07	0.05	0.10	0.08
dst_host_srv_diff_host_rate	0.00	0.12	0.00	0.00	0.00
dst_host_serror_rate	0.00	1.00	0.00	0.01	0.00
dst_host_srv_serror_rate	0.00	0.94	0.00	0.00	0.00
dst_host_rerror_rate	0.05	0.00	0.00	0.36	0.00
dst_host_srv_rerror_rate	0.05	0.00	0.00	1.00	0.00
attack	back	neptune	normal	portsweep	smurf

ตารางที่ 6.13 แสดงผลการแปลงค่าคุณลักษณะที่มีความต่อเนื่องของข้อมูลที่ใช้สอนและทดสอบระบบแต่ละชุดให้มีค่าอยู่ระหว่าง 0 ถึง 1 โดยคุณลักษณะที่ถูกแปลงค่าประกอบด้วยคุณลักษณะ duration count และ same_srv_rate เป็นต้น

2. นำคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์จากข้อ 1. ไปแปลงค่าเป็นตัวเลข หลังจากทำการแปลงค่าจะได้ผลดังตารางที่ 6.14

ตารางที่ 6.14 แสดงตัวอย่างข้อมูลหลังแปลงคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์ให้ตัวเลข

คุณลักษณะ	ตัวอย่างข้อมูล				
	duration	0.00	0.00	0.00	0.00
protocol_type	0	0	0	0	2
service	0	11	0	11	9
flag	0	4	0	7	0
src_bytes	0.00	0.00	0.00	0.00	0.00
dst_bytes	0.00	0.00	0.00	0.00	0.00
Land	0.00	0.00	0.00	0.00	0.00

ตารางที่ 6.14 (ต่อ)

คุณลักษณะ	ตัวอย่างข้อมูล				
wrong_fragment	0.00	0.00	0.00	0.00	0.00
urgent	0.00	0.00	0.00	0.00	0.00
count	0.00	0.03	0.01	0.03	1.00
srv_count	0.00	0.03	0.01	0.00	1.00
serror_rate	0.00	0.94	0.00	0.00	0.00
srv_serror_rate	0.00	1.00	0.00	0.00	0.00
rerror_rate	0.00	0.00	0.00	1.00	0.00
srv_rerror_rate	0.00	0.00	0.00	1.00	0.00
same_srv_rate	1.00	0.94	1.00	0.11	1.00
diff_srv_rate	0.00	0.12	0.00	0.50	0.00
srv_diff_host_rate	0.00	0.00	0.00	0.00	0.00
dst_host_count	1.00	0.06	0.07	0.74	0.62
dst_host_srv_count	1.00	0.06	0.07	0.00	0.05
dst_host_same_srv_rate	1.00	1.00	1.00	0.01	0.08
dst_host_diff_srv_rate	0.00	0.00	0.00	0.06	0.02
dst_host_same_src_port_rate	0.00	0.07	0.05	0.10	0.08
dst_host_srv_diff_host_rate	0.00	0.12	0.00	0.00	0.00
dst_host_serror_rate	0.00	1.00	0.00	0.01	0.00
dst_host_srv_serror_rate	0.00	0.94	0.00	0.00	0.00
dst_host_rerror_rate	0.05	0.00	0.00	0.36	0.00
dst_host_srv_rerror_rate	0.05	0.00	0.00	1.00	0.00
attack	back	neptune	normal	portsweep	smurf

ตารางที่ 6.14 แสดงผลการแปลงค่าคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์จากข้อ 1. เป็นตัวเลข คุณลักษณะที่มีลักษณะเป็นสัญลักษณ์ประกอบด้วย protocol feature service feature และ flag feature ตามลำดับ

3. นำข้อมูลที่ใช้สอนระบบแต่ละชุดจากข้อ 2. ไปสอนระบบตรวจจับการบุกรุกผ่านเครือข่าย โดยใช้นิรอลเน็ตเวิร์ค โดยในการสอนจะใช้วิธีการเรียนรู้แบบแพร่ย้อนกลับ (Back propagation) ซึ่งเป็นอัลกอริทึมที่ถูกออกแบบมาเพื่อใช้กับนิรอลเน็ตเวิร์คแบบหลายเลเยอร์และกำหนด

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับใช้งานเพื่อการศึกษาค้นคว้าเท่านั้น ไม่อนุญาตให้นำไปใช้เพื่อวัตถุประสงค์อื่น การค้า
ค่าพารามิเตอร์ในการทดลองดังตารางที่ 6.15 โดยจะทำการสอนนิรอลเน็ตเวิร์คจนกว่าจะได้ค่า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Mean Square Error น้อยกว่าหรือเท่ากับ 0.001 หรือครบ 2000 epochs หลังจากทำการสอนจะได้ค่า Mean Square Error ดังตารางที่ 6.16

ตารางที่ 6.15 พารามิเตอร์ในการสอนระบบเริ่มต้น

Parameter	Value
Network type	Feed-forward backpropagation
Training Function	Variable learning rate backpropagation (Traingdx)
Adaption learning function	Gradient descent with momentum (Learnngdm)
Performance function	Mean square error
Epochs	2000
Goal	0.001
Learning rate	0.01
Learning rate decerment	0.7
Learning rate increment	1.05
Maximum fail	5
Maximum performance increment	1.04
Momentum constant	0.95
Minimum grad	1e-006

ตารางที่ 6.16 แสดงค่า Mean Square Error ที่ทดลองได้จากข้อมูลแต่ละชุด

ชุดข้อมูล (สอน/ทดสอบ)	Mean Square Error (%)
20/80	0.00281276
40/60	0.00259635
50/50	0.00204127
60/40	0.00242398
80/20	0.00178475

ตารางที่ 6.16 แสดงผลการสอนนิรอลเน็ตเวิร์คด้วยข้อมูลที่ใช้สอนระบบแต่ละชุด จากตารางจะพบว่านิรอลเน็ตเวิร์คมีค่าเปอร์เซ็นต์ในการเรียนรู้ต่ำ เนื่องจากมีค่า Mean Square Error ที่ได้จากการสอนด้วยข้อมูลที่ใช้สอนระบบแต่ละชุดสูง

4. นำข้อมูลที่ใช้ทดสอบระบบแต่ละชุดไปทดสอบระบบ โดยผลการทดลองแสดงในตารางที่

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่นอญูตให้นำไปใช้ประโยชน์ด้านการค้า 6.17 และ 6.18 ตามลำดับ

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 6.17 Detection rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิวรอลเน็ตเวิร์ค

Attack Type	Detection Rate (%)					Average
	20/80	40/60	50/50	60/40	80/20	
back	3.64	10.30	0.00	3.24	0.00	3.44
ipsweep	92.35	0.00	0.00	0.00	94.40	37.35
land	0.00	0.00	0.00	0.00	0.00	0.00
neptune	99.99	100.00	99.99	100.00	100.00	100.00
nmap	45.65	0.00	0.00	47.83	0.00	18.70
normal	99.94	99.83	99.94	99.94	99.97	99.92
pod	0.00	0.00	0.00	0.00	96.30	19.26
portsweep	0.00	0.00	74.05	98.58	0.00	34.53
satan	0.00	94.44	94.96	88.21	87.42	73.01
smurf	99.98	99.99	99.99	99.83	100.00	99.96
teardrop	99.49	0.00	100.00	100.00	100.00	79.90
Average	49.19	36.78	51.72	57.97	61.64	51.46

ตารางที่ 6.17 แสดงผลการวัดค่า Detection rate ของระบบตรวจจับการบุกรุกโดยใช้นิวรอลเน็ตเวิร์คในการตรวจจับพฤติกรรมปกติและพฤติกรรมกรบุกรุกจำนวน 11 ประเภท จากตารางจะพบว่าระบบมีเปอร์เซ็นต์ความแม่นยำในการตรวจจับพฤติกรรมกรบุกรุก back, ipsweep, land, nmap, pod, และ portsweep ต่ำ เนื่องจากระบบไม่สามารถเรียนรู้และจำแนกพฤติกรรมกรบุกรุกประเภทดังกล่าวออกจากพฤติกรรมกรบุกรุกประเภทอื่นได้ ซึ่งจากตารางที่ 6.16 จะพบว่าระบบมีค่า Mean Square Error สูงทำให้มีเปอร์เซ็นต์ความแม่นยำในการตรวจจับพฤติกรรมกรบุกรุกประเภทดังกล่าวต่ำ โดยพฤติกรรมกรบุกรุกแต่ละประเภทจะมีค่า Detection rate เฉลี่ยต่ำกว่า 50.00% เมื่อนำค่า Detection rate ในการตรวจจับพฤติกรรมกรบุกรุกของข้อมูลแต่ละชุดมาคำนวณหาค่าเฉลี่ยจะได้ค่า Detection rate เฉลี่ยเท่ากับ 51.46% ดังนั้นจึงสรุปได้ว่าระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิวรอลเน็ตเวิร์คเป็นระบบที่มีเปอร์เซ็นต์ความแม่นยำในการตรวจจับพฤติกรรมกรบุกรุกต่ำ

ตารางที่ 6.18 False negative rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิวรอลเน็ตเวิร์ค

ชุดข้อมูล (สอน/ทดสอบ)	False negative rate (%)
20/80	3.489

ตารางที่ 6.18 (ต่อ)

ชุดข้อมูล (สอน/ทดสอบ)	False negative rate (%)
40/60	2.946
50/50	3.693
60/40	3.221
80/20	2.791
Average	3.228

ตารางที่ 6.18 แสดงผลการวัดค่า False negative rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิวรอลเน็ตเวิร์ค ซึ่งวัดจากการตรวจจับพฤติกรรมที่ผิดพลาดของของกฎที่ใช้จำแนกพฤติกรรมปกติ จากตารางจะพบว่าระบบมีค่า False negative rate ของข้อมูลแต่ละชุดมีสูงกว่า 2.000% เมื่อนำค่า False negative rate ของข้อมูลแต่ละชุดมาคำนวณหาค่าเฉลี่ยจะได้ค่า False negative rate เฉลี่ยเท่ากับ 3.228% ดังนั้นจึงสรุปได้ว่าระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้นิวรอลเน็ตเวิร์คเป็นระบบที่มีเปอร์เซ็นต์ความผิดพลาดในการตรวจจับพฤติกรรมบุกรุกสูง

6.4 การตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ดีซีชันทรี

ดีซีชันทรี คือแบบจำลองที่มีโครงสร้างคล้ายกับต้นไม้ที่จำแนกประเภทข้อมูลตามค่าของคุณลักษณะ โดยจะแทนกฎที่จำแนกประเภทข้อมูลโดยใช้ค่าของคุณลักษณะ ดีซีชันทรีประกอบด้วยโหนด (nodes) ลีฟ (leaves) และเส้นเชื่อมต่อระหว่างโหนด (edges) โหนดจะระบุคุณลักษณะโดยค่าของข้อมูลที่ถูกแบ่ง โดยแต่ละโหนดจะมีเส้นเชื่อมต่อระหว่างโหนดที่ถูกเลเบลด้วยค่าของคุณลักษณะที่อยู่ในพารินท์โหนด (parent node) เส้นเชื่อมต่อระหว่างโหนดจะเชื่อมโหนดและลีฟเข้าด้วยกัน ลีฟถูกเลเบลด้วยค่าตัดสินใจสำหรับจำแนกประเภทของข้อมูล ดีซีชันทรีนิยมนำไปใช้แก้ปัญหาเกี่ยวกับการจำแนกประเภทข้อมูล

การใช้ดีซีชันทรีในงานวิจัยนี้มีจุดประสงค์เพื่อเปรียบเทียบประสิทธิภาพในการจำแนกพฤติกรรมการบุกรุกกับระบบตรวจจับการบุกรุกที่นำเสนอในงานวิจัยนี้ สำหรับการทดลองจะใช้ข้อมูลชุดเดียวกันกับที่ใช้กับระบบที่นำเสนอในงานวิจัยนี้ (ข้อมูลในข้อ 6.1.2) ซึ่งได้มีการเตรียมข้อมูลสำหรับใช้สอนและทดสอบระบบ 5 ชุดดังตารางที่ 6.3 โดยจะทำการทดสอบระบบกับข้อมูลทั้ง 5 ชุดและนำค่า Detection rate และ False negative rate ของข้อมูลแต่ละชุดมาหาค่าเฉลี่ย โดยมีขั้นตอนการทดลองดังนี้

1. นำคุณลักษณะที่มีความต่อเนื่องของข้อมูลที่ใช้สอนและทดสอบระบบแต่ละชุดไป

แปลงให้มีค่าอยู่ระหว่าง 0 ถึง 1 โดยใช้สมการ 5.2 ในบทที่ 5 หัวข้อ 5.3.2

2. นำคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์จากข้อ 1. ไปแปลงค่าเป็นตัวเลข

3. นำข้อมูลที่ใช้สอนระบบแต่ละชุดไปสอนระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้
นิวรอลเน็ตเวิร์ค

4. นำข้อมูลที่ใช้ทดสอบระบบแต่ละชุดไปทดสอบระบบ

ต่อไปจะอธิบายรายละเอียดในแต่ละขั้นตอน

1. นำคุณลักษณะที่มีความต่อเนื่องของข้อมูลที่ใช้สอนและทดสอบระบบแต่ละชุดไปแปลง
ให้มีค่าอยู่ระหว่าง 0 ถึง 1 โดยใช้สมการ 5.2 ในบทที่ 5 หัวข้อ 5.3.2 ซึ่งหลังจากทำการแปลงค่าจะ
ได้ผลดังตารางที่ 6.19

ตารางที่ 6.19 แสดงตัวอย่างข้อมูลที่ถูกแปลงค่าคุณลักษณะที่มีความต่อเนื่องให้มีค่าอยู่ระหว่าง 0
ถึง 1

คุณลักษณะ	ตัวอย่างข้อมูล				
	0.00	0.00	0.00	0.00	0.00
duration	0.00	0.00	0.00	0.00	0.00
protocol_type	tcp	tcp	tcp	tcp	icmp
service	http	private	http	private	ecr_i
flag	SF	S0	SF	RSTR	SF
src_bytes	0.00	0.00	0.00	0.00	0.00
dst_bytes	0.00	0.00	0.00	0.00	0.00
Land	0.00	0.00	0.00	0.00	0.00
wrong_fragment	0.00	0.00	0.00	0.00	0.00
urgent	0.00	0.00	0.00	0.00	0.00
count	0.00	0.03	0.01	0.03	1.00
srv_count	0.00	0.03	0.01	0.00	1.00
serror_rate	0.00	0.94	0.00	0.00	0.00
srv_serror_rate	0.00	1.00	0.00	0.00	0.00
rerror_rate	0.00	0.00	0.00	1.00	0.00
srv_rerror_rate	0.00	0.00	0.00	1.00	0.00
same_srv_rate	1.00	0.94	1.00	0.11	1.00
diff_srv_rate	0.00	0.12	0.00	0.50	0.00
srv_diff_host_rate	0.00	0.00	0.00	0.00	0.00
dst_host_count	1.00	0.06	0.07	0.74	0.62
dst_host_srv_count	1.00	0.06	0.07	0.00	0.05
dst_host_same_srv_rate	1.00	1.00	1.00	0.01	0.08

ตารางที่ 6.19 (ต่อ)

คุณลักษณะ	ตัวอย่างข้อมูล				
	back	neptune	normal	portsweep	smurf
dst_host_diff_srv_rate	0.00	0.00	0.00	0.06	0.02
dst_host_same_src_port_rate	0.00	0.07	0.05	0.10	0.08
dst_host_srv_diff_host_rate	0.00	0.12	0.00	0.00	0.00
dst_host_serror_rate	0.00	1.00	0.00	0.01	0.00
dst_host_srv_serror_rate	0.00	0.94	0.00	0.00	0.00
dst_host_rerror_rate	0.05	0.00	0.00	0.36	0.00
dst_host_srv_rerror_rate	0.05	0.00	0.00	1.00	0.00
attack	back	neptune	normal	portsweep	smurf

ตารางที่ 6.19 แสดงผลการแปลงค่าคุณลักษณะที่มีความต่อเนื่องของข้อมูลที่ใช้สอนและทดสอบระบบแต่ละชุดให้มีค่าอยู่ระหว่าง 0 ถึง 1 โดยคุณลักษณะที่ถูกแปลงค่าประกอบด้วยคุณลักษณะ duration count และ same_srv_rate เป็นต้น

2. นำคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์จากข้อ 1. ไปแปลงค่าเป็นตัวเลข หลังจากทำการแปลงค่าจะได้ผลดังตารางที่ 6.20

ตารางที่ 6.20 แสดงตัวอย่างข้อมูลหลังแปลงคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์ให้ตัวเลข

คุณลักษณะ	ตัวอย่างข้อมูล				
	back	neptune	normal	portsweep	smurf
duration	0.00	0.00	0.00	0.00	0.00
protocol_type	0	0	0	0	2
service	0	11	0	11	9
flag	0	4	0	7	0
src_bytes	0.00	0.00	0.00	0.00	0.00
dst_bytes	0.00	0.00	0.00	0.00	0.00
Land	0.00	0.00	0.00	0.00	0.00
wrong_fragment	0.00	0.00	0.00	0.00	0.00
urgent	0.00	0.00	0.00	0.00	0.00
count	0.00	0.03	0.01	0.03	1.00
srv_count	0.00	0.03	0.01	0.00	1.00
serror_rate	0.00	0.94	0.00	0.00	0.00
srv_serror_rate	0.00	1.00	0.00	0.00	0.00

ตารางที่ 6.20 (ต่อ)

คุณลักษณะ	ตัวอย่างข้อมูล				
rerror_rate	0.00	0.00	0.00	1.00	0.00
srv_error_rate	0.00	0.00	0.00	1.00	0.00
same_srv_rate	1.00	0.94	1.00	0.11	1.00
diff_srv_rate	0.00	0.12	0.00	0.50	0.00
srv_diff_host_rate	0.00	0.00	0.00	0.00	0.00
dst_host_count	1.00	0.06	0.07	0.74	0.62
dst_host_srv_count	1.00	0.06	0.07	0.00	0.05
dst_host_same_srv_rate	1.00	1.00	1.00	0.01	0.08
dst_host_diff_srv_rate	0.00	0.00	0.00	0.06	0.02
dst_host_same_src_port_rate	0.00	0.07	0.05	0.10	0.08
dst_host_srv_diff_host_rate	0.00	0.12	0.00	0.00	0.00
dst_host_serror_rate	0.00	1.00	0.00	0.01	0.00
dst_host_srv_serror_rate	0.00	0.94	0.00	0.00	0.00
dst_host_rerror_rate	0.05	0.00	0.00	0.36	0.00
dst_host_srv_rerror_rate	0.05	0.00	0.00	1.00	0.00
attack	back	neptune	normal	portsweep	smurf

ตารางที่ 6.20 แสดงผลการแปลงค่าคุณลักษณะที่มีลักษณะเป็นสัญลักษณ์จากข้อ 1. เป็นตัวเลข คุณลักษณะที่มีลักษณะเป็นสัญลักษณ์ประกอบด้วย protocol feature service feature และ flag feature ตามลำดับ

3. นำข้อมูลที่ใช้สอนระบบแต่ละชุดจากข้อ 2. ไปสอนระบบตรวจจับการบุกรุกผ่านเครือข่าย โดยใช้ดัชนีชี้วัดที่อยู่ที่อยู่ในโปรแกรม MATLAB เวอร์ชัน 6.5 และกำหนดค่าพารามิเตอร์ในการทดลองดังตารางที่ 6.21 หลังจากทำการสอนจะได้จำนวนโหนด คุณลักษณะและกฎที่ใช้จำแนกพฤติกรรมการบุกรุกดังตารางที่ 6.22

ตารางที่ 6.21 พารามิเตอร์ในการสร้างดัชนีชี้วัด

Parameter	Value
Method	classification
Splitmin	10
Prune	On

ตารางที่ 6.22 แสดงจำนวน โหนด คุณลักษณะและกฎที่ใช้จำแนกพฤติกรรมการบุกรุกที่ได้จากการทดลอง

ชุดข้อมูล (สอน/ทดสอบ)	จำนวนโหนด	จำนวนคุณลักษณะ	จำนวนกฎที่ใช้จำแนก พฤติกรรมการบุกรุก
20/80	161	23	81
40/60	209	23	105
50/50	261	23	131
60/40	263	25	132
80/20	291	24	146

ตารางที่ 6.22 แสดงผลการสอนระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ดีซีชันทรีด้วยข้อมูลที่ใส่สอนระบบแต่ละชุดข้อมูล จากตารางจำนวนของโหนดและคุณลักษณะจะขึ้นกับความแตกต่างของข้อมูล ส่วนจำนวนของกฎที่ใช้จำแนกพฤติกรรมการบุกรุกจะขึ้นอยู่กับจำนวนของโหนดที่อยู่ในดีซีชันทรี โดยกฎที่ใช้จำแนกพฤติกรรมการบุกรุกแสดงดังตารางที่ 6.23

ตารางที่ 6.23 แสดงตัวอย่างกฎที่ใช้จำแนกการบุกรุกที่ได้จากการเรียนรู้ของระบบ

Condition	Action
(same_srv<0.48) & (dst_host_diff_srv<0.105)	neptune
(same_srv<0.48) & (dst_host_diff_srv>0.105) & (count<0.02)	normal
(same_srv>0.48) & (protocol>1.5) & (count<0.025) & (dst_host_count<0.045)	ipsweep
(same_srv>0.48) & (protocol>1.5) & (count<0.025) & (dst_host_count>0.045)	normal
(same_srv>0.48) & (protocol>1.5) & (count>0.025)	smurf
...	...

4. นำข้อมูลที่ใช้ทดสอบระบบแต่ละชุดไปทดสอบระบบ โดยผลการทดลองแสดงในตารางที่ 6.24 และ 6.25 ตามลำดับ

ตารางที่ 6.24 Detection rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ดีซีชันทรี

Attack Type	Detection Rate (%)					Average
	20/80	40/60	50/50	60/40	80/20	
back	96.49	92.55	93.15	92.48	89.90	92.91
ipsweep	99.44	99.50	98.36	97.76	98.88	98.79
land	75.00	66.67	90.00	87.50	75.00	78.83

ตารางที่ 6.24 (ต่อ)

Attack Type	Detection Rate (%)					Average
	20/80	40/60	50/50	60/40	80/20	
neptune	99.99	99.99	99.99	99.98	99.99	99.99
nmap	94.02	94.20	99.13	96.74	97.83	96.38
normal	99.83	99.78	99.82	99.82	99.85	99.82
pod	97.62	97.47	96.21	100.00	100.00	98.26
portsweep	96.57	83.89	97.73	99.53	100.00	95.54
satan	95.83	97.59	98.49	97.80	97.17	97.38
smurf	99.99	99.99	99.99	99.99	100.00	99.99
teardrop	100.00	99.83	100.00	100.00	100.00	99.97
Average	95.89	93.77	97.53	97.42	96.24	96.17

ตารางที่ 6.24 แสดงผลการวัดค่า Detection rate ของระบบตรวจจับการบุกรุกผ่านเครือข่าย โดยใช้ดัชนีชี้วัดที่ตรวจจับพฤติกรรมปกติและพฤติกรรมการบุกรุกจำนวน 11 ประเภท จากตารางจะพบว่าระบบมีเปอร์เซ็นต์ความแม่นยำในการตรวจจับพฤติกรรมการบุกรุก neptune และ smurf สูงที่สุด โดยพฤติกรรมการบุกรุกทั้งสองมีค่า Detection rate เฉลี่ยเท่ากับ 99.99% และมีค่าเปอร์เซ็นต์ความแม่นยำในการตรวจจับพฤติกรรมการบุกรุก land ต่ำที่สุดโดยมีค่า Detection rate เฉลี่ยเท่ากับ 78.83% เมื่อนำค่า Detection rate ของข้อมูลแต่ละชุดมาคำนวณหาค่าเฉลี่ยจะได้ค่า Detection rate เฉลี่ยเท่ากับ 99.37% ดังนั้นจึงสรุปได้ว่าระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ดัชนีชี้วัดเป็นระบบที่มีเปอร์เซ็นต์ความแม่นยำในการตรวจจับพฤติกรรมปกติและพฤติกรรมการบุกรุกสูงกว่าระบบที่ใช้นิรอลเน็ตเวิร์คแต่ก็ยังคงต่ำกว่าระบบที่ใช้วิธีราฟฟิชซี่

ตารางที่ 6.25 False negative rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ดัชนีชี้วัด

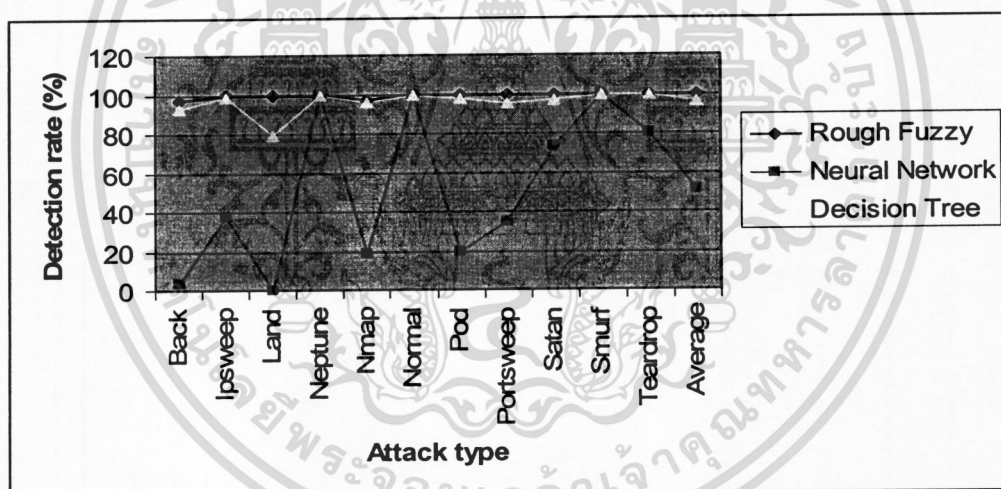
ชุดข้อมูล (สอน/ทดสอบ)	False negative rate (%)
20/80	0.203
40/60	0.246
50/50	0.238
60/40	0.240
80/20	0.289
Average	0.243

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 6.25 แสดงผลการวัดค่า False negative rate ของระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ดัชนีชี้วัด ซึ่งวัดจากการตรวจจับพฤติกรรมที่ผิดพลาดของของกฎที่ใช้จำแนกพฤติกรรมปกติ จากตารางจะพบว่าระบบมีค่า False negative rate เฉลี่ยเท่ากับ 0.243% ดังนั้นจึงสรุปได้ว่าระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ดัชนีชี้วัดเป็นระบบที่มีค่าเปอร์เซ็นต์ความผิดพลาดในการตรวจจับพฤติกรรมการบุกรุกต่ำแต่ก็ยังสูงกว่าระบบที่ใช้วิธีกราฟฟิชซี

6.5 เปรียบเทียบผลการตรวจจับของแต่ละระบบ

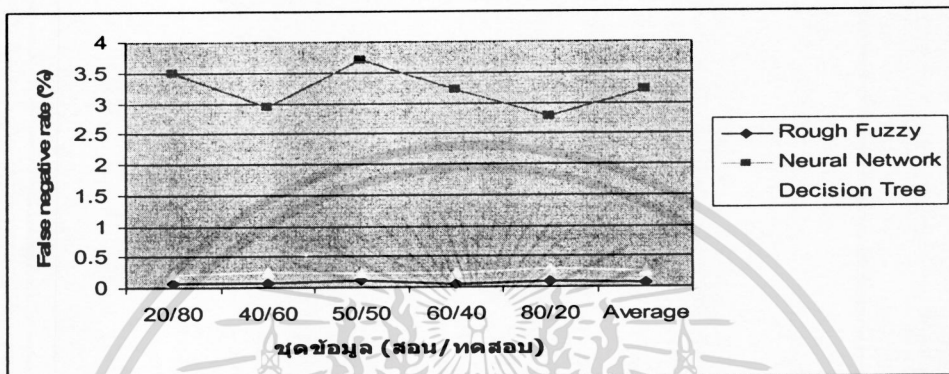
ในหัวข้อนี้จะเปรียบเทียบผลการทดสอบระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้กราฟฟิชซีที่นำเสนอในงานวิจัยนี้กับระบบตรวจจับการบุกรุกผ่านทางเครือข่ายโดยใช้นิวรอลเน็ตเวิร์คและดัชนีชี้วัด [13] ตามลำดับ โดยจะเปรียบเทียบผลการวัดค่า Detection rate และค่า False negative rate ของแต่ละระบบ ซึ่งผลการวัดค่า Detection rate และค่า False negative rate ของแต่ละระบบแสดงในรูปที่ 6.2 และ 6.3 ตามลำดับ



รูปที่ 6.2 กราฟแสดงค่า Detection rate เฉลี่ยของระบบตรวจจับการบุกรุกต่างๆ

รูปที่ 6.2 แสดงกราฟผลการวัดค่า Detection rate เฉลี่ยของระบบตรวจจับการบุกรุกแต่ละระบบ จากกราฟจะพบว่าระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้กราฟฟิชซีซึ่งมีค่า Detection rate เฉลี่ยในการตรวจจับพฤติกรรมปกติและพฤติกรรมการบุกรุกแต่ละประเภทสูงที่สุด โดยจะมีค่า Detection rate เฉลี่ยเท่ากับ 99.37% ส่วนระบบที่มีค่า Detection rate เฉลี่ยรองลงมาคือระบบที่ใช้ดัชนีชี้วัดมีค่า Detection rate เฉลี่ยเท่ากับ 96.17% โดยพฤติกรรมการบุกรุกที่ระบบที่ใช้ดัชนีชี้วัดวัดค่า Detection rate เฉลี่ยต่ำที่สุดคือพฤติกรรมการบุกรุกแบบ land มีค่าเท่ากับ 78.83% และระบบที่มีค่า Detection rate เฉลี่ยสูงที่สุดคือระบบที่ใช้นิวรอลเน็ตเวิร์ค โดยจะมีค่า Detection rate เฉลี่ยเท่ากับ 51.46% โดยพฤติกรรมการบุกรุกที่ระบบที่ใช้นิวรอลเน็ตเวิร์ควัดค่า Detection

rate เกลี่ยได้ต่ำคือพฤติกรรมการบุกรุกแบบ back, ipsweep, land, nmap, pod, และ portsweep โดยมีค่า Detection rate เกลี่ยเท่ากับ 3.44%, 37.35%, 0.00%, 18.70%, 19.26%, และ 34.53% ตามลำดับ ดังนั้นจึงสรุปได้ว่าระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซึ่งจะมีค่าเปอร์เซ็นต์ความแม่นยำในการตรวจจับพฤติกรรมปกติและพฤติกรรมการบุกรุกสูงกว่าระบบที่ใช้ดีซิชั่นทรีและระบบที่ใช้นิวรอลเน็ตเวิร์คตามลำดับ



รูปที่ 6.3 กราฟแสดงค่า False negative rate เกลี่ยของระบบตรวจจับการบุกรุกต่างๆ

รูปที่ 6.3 แสดงกราฟผลการวัดค่า False negative rate เกลี่ยของระบบตรวจจับการบุกรุกแต่ละระบบที่ทดลองกับข้อมูลที่ใช้ทดสอบแต่ละชุด จากกราฟจะพบว่าระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซึ่งจะมีค่า False negative rate เกลี่ยต่ำที่สุดคือจะมีค่าเท่ากับ 0.079% ส่วนระบบที่ใช้ดีซิชั่นทรีและนิวรอลเน็ตเวิร์คจะมีค่า False negative rate เกลี่ยเท่ากับ 0.243% และ 3.228% ตามลำดับ ดังนั้นจึงสรุปได้ว่าระบบตรวจจับการบุกรุกผ่านเครือข่ายโดยใช้ราฟฟิซซึ่งจะมีค่าเปอร์เซ็นต์ความผิดพลาดในการตรวจจับพฤติกรรมการบุกรุกของกฎที่ใช้จำแนกพฤติกรรมปกติต่ำที่สุดและระบบที่ใช้นิวรอลเน็ตเวิร์คสูงที่สุด

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] R. Baeza-Yates and B. Ribeiro-Neto., "Modern Information Retrieval," New York: ACM-Press. 1999.
- [9] S. Stolfo et al., The Third International Knowledge Discovery and Data Mining Tools Competition., <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>
- [10] W. Lee and S. J. Stolfo, "A Framework for Constructing Features and Models for Intrusion Detection Systems", In ACM Transactions on Information and System Security, Vol. 3, No. 4, November 2000, pp. 227-261
- [11] M. Mohajerani, A. Moeini and M. Kianie, "NFIDS: A Neuro-fuzzy Intrusion Detection System", In Proceedings of the 10th IEEE International Conference on Electronics, Circuits and Systems, 2003, pp348-351.

- [12] J. Gomez and D. Dasgupta, "Evolving Fuzzy Classifiers for Intrusion Detection", In Proceedings of the 2002 IEEE Workshop on Information Assurance United States Military Academy, West Point, NY June 2001
- [13] Z.S. Pan , S.C.Chen , G.B.Hu , and D.Q.Zhang, "Hybrid Neural Network and C4.5 for Misuse Detection", In Proceedings of the second International Conference on Machine Learning and Cybernetics, Xi'an , 2-5 November 2003.
- [14] J.K. Liang, Y.Zhang, Y.B. Qu, "A Heuristic Algorithm of Attribute Reduction in Rough Set", In Proceedings of the Fourth International Conference on Machine Learning and Cybernetics, Guangzhou, 18-21 August 2005.
- [15] J. Huang, S. Li, C. Man, "A T-S Type of Rough Fuzzy Controller Based on Process Input-output Data", In Proceedings of the 42nd IEEE conference on Decision and control Maui, Hawaii USA, December 2003
- [16] Z. Lian-hua, Z Guan-hua, Y. Lang, "Intrusion detection using rough set classification" ,In Journal of Zhejiang University SCIENCE, pp.1076-1086, 2004.
- [17] J.Komorowski , Z.Pawlak , L. Polkowski , and A.Skowron , "Rough Sets: A Tutorial" , Rough-Fuzzy Hybridization: A New Trend in Decision Making , pp. 3-98 , 1999
- [18] George J. Klir et. Al. , Fuzzy Set Theory:Foundations and Applications , Prentice Hall International , Inc.



ภาคผนวก

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



วิศวกรรมลาดกระบัง

Ladkrabang Engineering Journal

คณะวิศวกรรมศาสตร์ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง กรุงเทพฯ 10520
Faculty of Engineering, King Mongkut's Institute of Technology Ladkrabang, Bangkok 10520

วันที่ 4 มกราคม 2550

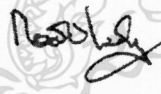
เลขที่อ้างอิง 1161

เรื่อง การตอบรับบทความ

เรียน คุณธรรมกร ครองไตรภพ เอื้อน ปิ่นเงิน

ตามที่ท่านได้ส่งบทความเรื่อง การตรวจจับการบุกรุกโดยใช้ราฟฟิซซี่ (Rough Fuzzy Intrusion Detection System) มาให้พิจารณาเพื่อลงตีพิมพ์ในวารสารวิศวกรรมลาดกระบัง บัดนี้ ผู้ทรงคุณวุฒิได้ทำการพิจารณาแล้วเห็นว่า ยอมรับตีพิมพ์ได้ โดยจะตีพิมพ์ในปีที่ 23 ฉบับที่ 4 เดือน ธันวาคม 2549

จึงเรียนมาเพื่อทราบ


(รศ.ดร.กอบชัย เดชธานี)
หัวหน้ากองบรรณาธิการ

การขยายราฟเซตโดยใช้ฟัซซีเซตสำหรับระบบตรวจจับการบุกรุก

Extended Rough Set Classification Using Fuzzy Set for Intrusion Detection System

ธรรมกร ครองไทรภพ เอื้อน ปิ่นเงิน

คณะวิศวกรรมศาสตร์และสาขาคอมพิวเตอร์สารสนเทศ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

บทคัดย่อ

ปัจจุบันระบบตรวจจับการบุกรุกเครือข่ายโดยใช้เทคนิคการเรียนรู้ของเครื่องจักร (machine learning) เป็นหัวข้องานวิจัยที่แพร่หลาย เนื่องจากสามารถตรวจจับได้ทั้งการบุกรุกแบบ misuse และ anomaly งานวิจัยนี้ได้เสนอระบบตรวจจับการบุกรุกโดยใช้ราฟเซต ซึ่งต่างกับวิธีราฟเซตเดิมตรงที่วิธีราฟเซตจะทำการแบ่งช่วงลักษณะเด่นที่มีความต่อเนื่องโดยใช้วิธีแบ่งความกว้างช่วงค่าที่เท่ากันส่วนวิธีที่เสนอนี้ใช้ฟัซซีเซตและราฟเซตจัดการกับความคลุมเครือและความไม่แน่นอนของข้อมูล ระบบที่เสนอใช้ราฟเซตในการสร้างกฎที่ใช้จำแนกการบุกรุกสำหรับจำแนกพฤติกรรมบุกรุกและพฤติกรรมปกติ โดยแบ่งการทำงานออกเป็นสี่ขั้นตอนหลัก ขั้นตอนแรกเป็นการสร้างระบบตัดสินใจโดยใช้ข้อมูลจากฐานข้อมูลของ KDD ขั้นตอนที่สองแบ่งช่วงลักษณะเด่นที่มีความต่อเนื่องของข้อมูลโดยใช้ฟัซซีเซต ขั้นตอนที่สามนำข้อมูลที่ได้มากำจัดคุณลักษณะที่ไม่มีความจำเป็นในการจำแนกและสร้างกฎที่ใช้จำแนกการบุกรุกสำหรับจำแนกพฤติกรรมบุกรุกและพฤติกรรมปกติ และขั้นตอนสุดท้ายทดสอบประสิทธิภาพของกฎที่ใช้จำแนกการบุกรุก จากการทดลองระบบด้วยฐานข้อมูล KDD 99 พบว่า Detection rate เฉลี่ยเท่ากับ 99.14% และ False negative rate เท่ากับ 0.022% ซึ่งมีความแม่นยำสูงกว่าระบบที่ใช้ราฟเซต

Abstract

Recently machine learning-based intrusion detection approaches have been subjected to extensive researches because they can detect both misuse and anomaly. In this paper, we present an intrusion detection system based on rough-fuzzy approach. Unlike other existing rough set approaches, rough set is used to discretize continuous data to subintervals of equal width. For proposed system, we use rough set to create an intrusion rules for classifying normal and abnormal behaviors. The process consists of four stages. First, create the decision system using data from KDD database. Second, discretize the continuous data to subintervals using fuzzy set. Third, reduce the redundancy, inconsistency of initial information table and generate intrusion rules using rough set. Finally, test the efficiency of intrusion rules using KDD 99 data set. From the experiments, the average of the detection rate is 99.14% and the false negative rate 0.022% higher accuracy than rough set.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1. บทนำ

ปัจจุบันการตรวจจับการบุกรุกมีความจำเป็นอย่างมาก เนื่องจากการบุกรุกระบบคอมพิวเตอร์เพิ่มมากขึ้นตามจำนวนของคอมพิวเตอร์ที่ต่อเข้าสู่เครือข่าย ซึ่งก่อให้เกิดความเสียหายกับองค์กรทางธุรกิจต่างๆ ดังนั้นความสนใจในการพัฒนาระบบตรวจจับการบุกรุกเครือข่ายจึงมุ่งไปที่การใช้เทคนิคการเรียนรู้ของเครื่องจักร ซึ่งเป็นหัวข้องานวิจัยที่แพร่หลาย การตรวจจับการบุกรุก (Intrusion Detection: ID) [1] เป็นการตรวจจับบุคคลผู้ที่ใช้ระบบคอมพิวเตอร์โดยไม่ได้รับการอนุญาต และบุคคลผู้ที่สามารถเข้าไปในระบบตามสิทธิที่ได้รับ แต่ใช้สิทธิที่มีอยู่ในทางที่ผิด คำนิยามที่เหมาะสมควรเป็น “การตรวจจับความพยายามที่จะใช้ระบบคอมพิวเตอร์โดยไม่ได้รับอนุญาตหรือใช้สิทธิที่มีอยู่ในทางที่ผิด”

วิธีตรวจจับการบุกรุก แบ่งออกได้ 2 วิธีหลักๆ คือ misuse detection และ anomaly detection [1] วิธี misuse detection พยายามจะจำการบุกรุกที่ได้มาจากรูปแบบที่ได้รับรองและรายงานโดยผู้เชี่ยวชาญ วิธี misuse detection มีช่องโหว่ (vulnerable) ตรงที่ผู้บุกรุกจะใช้รูปแบบของพฤติกรรมใหม่หรือผู้ที่ปิดบังพฤติกรรมที่ไม่ถูกต้องเพื่อหลอกลวงระบบตรวจจับ วิธี anomaly detection ถูกพัฒนาเพื่อแก้ปัญหานี้ โดยวิธี anomaly detection จะแทนรูปแบบของพฤติกรรมปกติ โดยได้สมมติว่าการบุกรุกสามารถตรวจจับได้โดยใช้พฤติกรรมที่เบี่ยงเบนไปจากพฤติกรรมปกติ สัญญาณแจ้งเตือนการบุกรุกจะถูกสร้างขึ้น เมื่อพฤติกรรมที่เบี่ยงเบนถูกตรวจจับได้

เมื่อไม่นานมานี้ นักวิจัยได้พัฒนาวิธีการตรวจจับการบุกรุกขึ้นเป็นจำนวนมาก อาทิเช่น Lee และ Stolfo [2] ได้เสนอวิธีตรวจจับการบุกรุกโดยใช้เทคนิคค่าไคโมนิ่งคำนวณหารูปแบบพฤติกรรมจากข้อมูลที่ตรวจสอบจากระบบและดึงคุณลักษณะที่ใช้ทำนายการบุกรุกออกจากรูปแบบพฤติกรรม จากนั้นใช้เทคนิคการเรียนรู้ของเครื่องจักรสร้างกฎตามนิยามของคุณลักษณะ Gomez และ Dasgupta [3] ได้เสนอวิธีตรวจจับการบุกรุกโดยใช้เจเนติกอัลกอริทึมสร้างกฎพีชชีสำหรับตรวจจับการบุกรุกโดยใช้ชุดข้อมูลที่มีรูปแบบพฤติกรรมของระบบที่เป็น

เงื่อนไข normal และ abnormal Lee และ Heinbuch [4] ได้เสนอวิธีตรวจจับการบุกรุกโดยใช้ neural network แบบ hierarchical backpropagation ตรวจจับการโจมตีแบบ neptune และ portsweep งานวิจัยเหล่านี้ยังมี Detection rate ต่ำ และ False alarm สูง นอกจากนี้กฎที่สร้างยังไม่สามารถอธิบายได้

งานวิจัยนี้ ได้นำเสนอวิธีตรวจจับการบุกรุกแบบ misuse โดยใช้วิธีราฟพีชชี เนื่องจากระบบตรวจจับการบุกรุกส่วนมากถูกสร้างโดยใช้รูปแบบพฤติกรรมการบุกรุกที่ได้รับรองและรายงานโดยผู้เชี่ยวชาญ และรูปแบบพฤติกรรมการบุกรุกที่มีอยู่ในปัจจุบันเพียงพอที่จะใช้ตรวจจับการบุกรุกใหม่ๆ ได้ นอกจากนี้วิธีตรวจจับการบุกรุกที่ใช้วิธีราฟพีชชียังสามารถสร้างกฎที่ใช้จำแนกการบุกรุกที่มีความแน่นอนสามารถอธิบายได้และมี Detection rate สูง อันดับแรกจะใช้พีชชีแยกช่วงลักษณะเด่นที่มีความต่อเนื่องของข้อมูล จากนั้นจะใช้ราฟพีชชีกำจัดคุณลักษณะที่ไม่จำเป็นในการจำแนกโดยคำนวณหา reducts หรือ minimal subsets ของคุณลักษณะและสร้างกฎที่ใช้จำแนกการบุกรุก โดยใช้ discernibility matrix และ boolean reasoning และทำการทดลองกับการบุกรุกระบบเครือข่ายแบบ Denial of service กับ Probing และพฤติกรรมปกติ

เนื้อหาส่วนที่เหลือของบทความนี้ ประกอบด้วยเนื้อหาต่อไปนี้ หัวข้อที่ 2 ทฤษฎีพื้นฐานของราฟพีชชี หัวข้อที่ 3 ขั้นตอนการทำงานของระบบตรวจจับการบุกรุกโดยวิธีราฟพีชชี หัวข้อที่ 4 การทดลองและอธิบายผล หัวข้อที่ 5 บทสรุป หัวข้อที่ 6 กิจกรรมประกาศ หัวข้อที่ 7 เอกสารอ้างอิง

2. ทฤษฎีพื้นฐานของราฟพีชชี

ทฤษฎีราฟพีชชีถูกคิดค้นในช่วงต้นทศวรรษ 1980 โดย Zdzislaw Pawlak (Pawlak, 1982) เป็นเครื่องมือทางคณิตศาสตร์ที่ใช้กับระบบผู้ช่วยในการตัดสินใจ (decision support system) และเหมาะที่จะใช้ในการจำแนกวัตถุมาก นอกจากนี้ราฟพีชชียังสามารถใช้ในการคัดเลือกคุณลักษณะ (feature selection) ดึงคุณลักษณะ (feature extraction) ฯลฯ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

นิยามที่ 1 ระบบสารสนเทศ S นิยามได้ดังนี้

$$S = (U, A, V, f) \tag{1}$$

เมื่อ $U = \{x_1, \dots, x_n\}$ คือเซตของวัตถุ (n คือจำนวนของวัตถุ), $A = \{a_1, \dots, a_m\}$ เซตของคุณลักษณะ (m คือจำนวนของคุณลักษณะ), $A = C \cup D$ และ $C \cap D = \emptyset$, C และ D คือเซตของคุณลักษณะที่ใช้ในการสร้างเงื่อนไข (condition attribute) และคุณลักษณะที่ใช้ในการตัดสินใจ (decision attribute) ตามลำดับ, $V = \bigcup_{a \in A} V_a$ และ V_a คือเซตค่าของคุณลักษณะ a ที่อยู่ใน A , และ f คือฟังก์ชัน $f: U \times V \rightarrow V$

นิยามที่ 2 กำหนดให้ $S = (U, A, V, f)$ เป็นระบบสารสนเทศ discernibility matrix M นิยามได้ดังนี้

$$(m_{ij}) = \begin{cases} \{a \in A : a(x_i) \neq a(x_j)\}, & D(x_i) \neq D(x_j) \\ \emptyset, & D(x_i) = D(x_j) \end{cases} \tag{2}$$

เมื่อ $a(x)$ คือค่าของวัตถุ x ณ คุณลักษณะ a , $D(x)$ คือค่าของของวัตถุ x ณ decision attribute D และ (m_{ij}) คือสมาชิกของ M เมื่อ $i, j = 1, 2, 3, \dots, n$

นิยามที่ 3 กำหนดให้ $S = (U, A, V, f)$ เป็นระบบสารสนเทศ Discernibility function f_S นิยามได้ดังนี้

$$f_S(a_1, \dots, a_m) = \bigwedge \{ \bigvee c_{ij} \mid 1 \leq j \leq i \leq n, c_{ij} \neq \emptyset \} \tag{3}$$

เมื่อ $\bigvee (c_{ij})$ คือ disjunction ของตัวแปร a และ $a \in c_{ij}$.

นิยามที่ 4 กำหนดให้ $S = (U, A \cup \{d\})$ เป็นระบบสารสนเทศ Decision rule นิยามได้ดังนี้

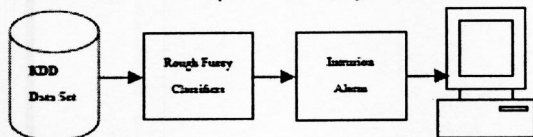
$$(a_1 = v_1) \wedge \dots \wedge (a_n = v_n) \Rightarrow (d = k) \tag{4}$$

เมื่อ $1 \leq i_1 < \dots < i_m \leq |A|, v_j \in V_{a_j}$ และ k คือค่าตัดสินใจ

3. ขั้นตอนการทำงานของระบบ

3.1 ขั้นตอนการทำงานของระบบ

ขั้นตอนการทำงานของระบบตรวจจัดการบุกรุกที่นำวิธีการฟัซซี่เข้ามาประยุกต์ใช้แสดงดังรูปที่ 1



รูปที่ 1 แสดงขั้นตอนการทำงานของระบบ

จากรูปส่วนแรก คือข้อมูลที่ใช้ในงานวิจัยนี้เป็นข้อมูลที่ได้จาก KDD cup[5] เป็นข้อมูลที่จำลองพฤติกรรมการบุกรุก

รุกราน 4 ประเภทหลักดังนี้

1. Denial of Services Attacks คือการบุกรุกที่ผู้บุกรุกพยายามทำให้ระบบหยุดให้บริการ เช่น Mailbomb, SYN Flood, และ Smurf เป็นต้น
2. Probing คือการบุกรุกที่ผู้บุกรุกพยายามตรวจสอบจุดอ่อนของระบบ เช่น Nmap, Saint, และ Satan เป็นต้น
3. Remote to User Attacks คือการบุกรุกที่ผู้บุกรุกไม่มีสิทธิ์อยู่ในระบบแต่พยายามทำให้ตัวเองเข้าสู่ระบบได้ เช่น Guest, Imap, และ sendmail เป็นต้น
4. User to Root Attacks คือการบุกรุกที่ผู้บุกรุกพยายามเข้าสู่ระบบโดยใช้สิทธิ์เท่าเทียมกับ root เช่น Eject, Fdformat, และ Loadmodule เป็นต้น

นอกจากนี้ยังมีคุณลักษณะอีก 41 คุณลักษณะได้มาจากการดักจับข้อมูลที่มีการสื่อสารกันในระบบเครือข่ายคอมพิวเตอร์แล้วใช้เทคนิคค่าไบนารีเปลี่ยนแปลงให้อยู่ในรูปแบบของตัวอักษร โดยแบ่งออกได้ 4 กลุ่มดังนี้ [5]

1. Basic features เป็นคุณลักษณะพื้นฐานที่ได้จากแพคเกจข้อมูลที่สื่อสารในเครือข่าย มีทั้งหมด 9 คุณลักษณะ เช่น เวลาในการเชื่อมต่อ เป็นต้น
2. Content features เป็นคุณลักษณะที่เก็บรวบรวมข้อมูลที่แสดงให้เห็นพฤติกรรมน่าสงสัย มีทั้งหมด 13 คุณลักษณะ เช่น ความผิดพลาดในการล็อกอิน เป็นต้น
3. Time-based features เป็นคุณลักษณะที่เก็บรวบรวมข้อมูลที่แสดงลักษณะของการสื่อสาร มีทั้งหมด 9 คุณลักษณะ เช่น จำนวนครั้งในการเชื่อมต่อเข้าสู่ระบบเมื่อผ่านไป 2 วินาที เป็นต้น
4. Host-based features เป็นคุณลักษณะที่เก็บรวบรวมข้อมูลที่แสดงลักษณะของการสื่อสารไปยังเครื่องปลายทาง เครื่องเดิมตลอดเวลา มีทั้งหมด 10 คุณลักษณะ เช่น จำนวนครั้งในการเชื่อมต่อไปยังเครื่องปลายทางเดิม เป็นต้น

ส่วนที่สองการตรวจจัดการบุกรุกโดยใช้วิธีการฟัซซี่เป็นส่วนที่นำวิธีการฟัซซี่มาประยุกต์ใช้ในการสร้างกฎสำหรับจำแนกพฤติกรรมบุกรุกและพฤติกรรมปกติแล้วส่งผลการจำแนกไปยังส่วนที่สามที่ทำหน้าที่แจ้งเตือนเพื่อดำเนินการแจ้งเตือนไปยังคอมพิวเตอร์

3.2 การตรวจจัดการบุกรุกโดยใช้วิธีการฟัซซี่

อัลกอริทึมของการตรวจจับการบุกรุกโดยใช้ราฟฟิซซี่

[8] มี 4 ขั้นตอน คือ

1. กระบวนการสร้างระบบคัดลินใจ
2. กระบวนการแปลงข้อมูลเบื้องต้น
3. กระบวนการสร้างกฎที่ใช้จำแนกการบุกรุก
4. กระบวนการทดสอบกฎที่ใช้จำแนกการบุกรุก

3.2.1 กระบวนการสร้างระบบคัดลินใจ

ในการสร้างระบบคัดลินใจ เพื่อใช้ในระบบตรวจจับการบุกรุกโดยวิธีราฟฟิซซี่นั้น ในงานวิจัยนี้ได้กำหนดบนพื้นฐานของประเภทการบุกรุกแบบ Denial of Service กับ Probing และพฤติกรรมปกติ กลุ่มข้อมูลมีจำนวนทั้งสิ้น 492,843 เรคคอร์ด และการจำแนกพฤติกรรมทั้งสามจะใช้คุณลักษณะเพียง 28 คุณลักษณะ [2] ประกอบด้วย 3 กลุ่มคือ Basic features , Time-based features , และ Host-based features โดยจะแบ่งข้อมูลออกเป็น 2 กลุ่มเพื่อใช้ในการสอนระบบจำนวน 394,274 เรคคอร์ด และใช้ในการทดสอบระบบจำนวน 98,569 เรคคอร์ดดังตารางที่ 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

3.2.2 กระบวนการแปลงข้อมูลเบื้องต้น

ในกระบวนการแปลงข้อมูลเบื้องต้น จะทำการแปลงชนิดข้อมูลของคุณลักษณะที่มีชนิดข้อมูลเป็น continuous ให้เป็น discrete หรือช่วง [3] โดยอันดับแรกจะทำการแปลงค่าของแต่ละคุณลักษณะให้มีค่าอยู่ระหว่าง 0 ถึง 1 โดยใช้สมการต่อไปนี้:

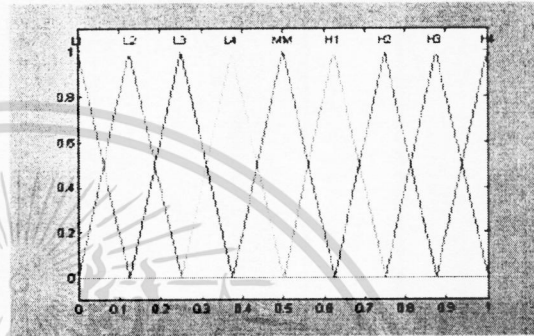
$$x = \frac{x - MIN}{MAX - MIN} \tag{5}$$

เมื่อ x คือค่าของคุณลักษณะ

MIN คือค่าต่ำที่สุดของคุณลักษณะ

MAX คือค่าสูงที่สุดของคุณลักษณะ

จากนั้นทำการแปลงค่าที่คำนวณได้เป็นค่าฟิซซี่ โดยใช้ fuzzy space ในรูปที่ 2



รูปที่ 2 Fuzzy space ของแต่ละ features

3.2.3 กระบวนการสร้างกฎที่ใช้จำแนกการบุกรุก

ในกระบวนการสร้างกฎที่ใช้จำแนกการบุกรุกนี้ จะทำการสร้างกฎที่ใช้จำแนกการบุกรุก โดยใช้ reducts หรือ minimal subsets ของคุณลักษณะที่คำนวณหาได้จากอัลกอริทึมลดคุณลักษณะ ในงานวิจัยนี้จะใช้อัลกอริทึม discernibility matrix และ Boolean calculation algorithm [6] [7] ซึ่งอัลกอริทึมจะมีขั้นตอนดังนี้

1. H1 discernibility matrix M
2. H1 discernibility function L จาก discernibility matrix M โดย discernibility function ที่ได้จะอยู่ในรูปแบบ disjunctive normal form (DNF)
3. ทำการลดรูป discernibility function L ที่ได้จากข้อ 2 โดยใช้ absorption laws
4. ทำการแปลง discernibility function L ที่ได้จากข้อ 3 ซึ่งอยู่ในรูปแบบ disjunctive normal form (DNF) ให้อยู่ในรูปแบบ conjunctive normal form (CNF) โดยใช้ distribution laws
5. ทำการลดรูป discernibility function L ที่ได้จากข้อ 4 โดยใช้ absorption laws
6. ทำการเลือกกลุ่มของ reduct ที่จะนำมาใช้สร้างกฎที่ใช้จำแนกการบุกรุก

เมื่อดำเนินการตามอัลกอริทึมเสร็จ จะได้กฎที่ใช้
จำแนกการบุกรุกของระบบตรวจจัดการบุกรุกโดยใช้ราฟ
ฟัซซี่ดังตารางที่ 2

Condition	Action
(duration=L1)&(wrong=L1)&(srv_count=L1L2)& (serror=H3H2)&(rerror=L1)&(srv_rerror=L1)& (same_srv=H3H2)&(srv_diff_host=L1)& (dst_host_count=L1L2)&(dst_host_srv_count= L1L2)&(dst_host_diff_srv=L1)& (dst_host_same_src_port=L3L2) &(service=private)& (dst_host_srv_diff=L4L3)&(dst_host_serror=H4)& (dst_host_srv_serror=H3H2)&(dst_host_rerror=L1)	neptune
(duration=L1)&(wrong=L1)&(srv_count=L1L2)& (serror=L1)&(rerror=L1)&(srv_rerror=L1)& (same_srv=H4)&(srv_diff_host=L1)& (dst_host_count=L1L2)&(dst_host_srv_count= L1L2)&(dst_host_diff_srv=L1)& (dst_host_same_src_port=L2L1)&(service=http)& (dst_host_srv_diff=L1)&(dst_host_serror=L1)& (dst_host_srv_serror=L1)&(dst_host_rerror=L1)	normal
...	...

ตารางที่ 2 ตัวอย่างกฎที่ใช้จำแนกการบุกรุก

3.2.4 การทดสอบกฎที่ใช้จำแนกการบุกรุก

ในการทดสอบกฎที่ใช้จำแนกการบุกรุกจะทำการ
ทดสอบกับกลุ่มข้อมูลที่ใช้ทดสอบ ซึ่งการทดสอบจะเน้น
เฉพาะค่า Detection Rate และ False Negative Rate เท่านั้น

Detection Rate คือค่าเปอร์เซ็นต์ของจำนวนการบุกรุก
ที่ตรวจจับ ได้ถูกต้องกับจำนวนของการบุกรุกที่ตรวจจับ ได้
ทั้งหมดที่ถูกระบุประเภทการบุกรุก โดยกฎที่ใช้จำแนกการ
บุกรุก จำนวนได้ดังนี้

$$\%detected = \left(\frac{N_{corrected}}{N_{corrected} + N_{missed}} \right) * 100 \quad (6)$$

เมื่อ $N_{corrected}$ คือจำนวนข้อมูลทดสอบที่กฎที่ใช้จำแนกการ
บุกรุกจำแนกได้ถูกต้อง

N_{missed} คือจำนวนข้อมูลทดสอบที่กฎที่ใช้จำแนกการ
บุกรุกจำแนกผิดพลาด

False Negative Rate คือค่าเปอร์เซ็นต์ของจำนวนการ
บุกรุกที่กฎที่ใช้จำแนกพฤติกรรมปกติจำแนกผิดพลาดกับ
จำนวนพฤติกรรมปกติที่จำแนกได้ถูกต้อง จำนวนได้ดังนี้

$$\%FalseNegative = \left(\frac{N_{false}}{N_{normal}} \right) * 100 \quad (7)$$

เมื่อ N_{false} คือจำนวนของการบุกรุกที่กฎที่ใช้จำแนก
พฤติกรรมปกติจำแนกผิดพลาด

N_{normal} คือจำนวนพฤติกรรมปกติที่กฎที่ใช้จำแนก
พฤติกรรมปกติจำแนกได้ถูกต้อง

4. ผลการทดลองและอธิบายผล

ตารางที่ 3 แสดงค่า Detection Rate ของระบบ
ตรวจจัดการบุกรุกในการตรวจจับพฤติกรรมการบุกรุก
จำนวน 11 ประเภทโดยรวมพฤติกรรมปกติด้วย จากการ
ทดลองพบว่าระบบตรวจจัดการบุกรุกโดยใช้ราฟฟัซซี่มีค่า
Detection Rate เฉลี่ยประมาณ 99.14% ซึ่งสูงกว่าระบบ
ตรวจจัดการบุกรุกโดยใช้ราฟเซต [8] ที่มีค่า Detection
Rate เฉลี่ยประมาณ 98.84%

Attack Type	Detection Rate	
	Rough Fuzzy	Rough set[8]
Back	99.15%	97.17%
Ipsweep	100.00%	100.00%
Land	100.00%	100.00%
Neptune	100.00%	100.00%
Nmap	92.59%	91.30%
Normal	99.46%	99.42%
Pod	100.00%	100.00%
Portswep	99.39%	99.35%
Satan	100.00%	100.00%
Smurf	100.00%	100.00%
Teardrop	100.00%	100.00%
Total	99.14%	98.84%

ตารางที่ 3 เปรียบเทียบ Detection Rate ของระบบ
ตรวจจัดการบุกรุกโดยใช้ราฟฟัซซี่กับราฟเซต

ตารางที่ 4 แสดงค่า False negative rate ซึ่งได้จากการ
ตรวจจับพฤติกรรมที่ผิดพลาดของระบบที่ตรวจจับ
พฤติกรรมปกติ จากการทดลองพบว่าระบบตรวจจัดการ
บุกรุกโดยใช้ราฟฟัซซี่จะมีค่า False negative rate ประมาณ
0.022% ซึ่งต่ำกว่าระบบตรวจจัดการบุกรุกโดยใช้ราฟเซต
[8] ที่มีค่า False negative rate ประมาณ 0.024%

IDS Type	False negative rate
Rough Fuzzy	0.022%
Rough set [8]	0.024%

ตารางที่ 4 False negative rate ของระบบตรวจจับการบุกรุกโดยใช้ราฟฟิชซีและราฟเซต

5. บทสรุป

งานวิจัยนี้ ได้เสนอระบบตรวจจับการบุกรุกโดยใช้ราฟฟิชซี ทฤษฎีฟัชซีเซตถูกใช้ในการแบ่งช่วงลักษณะเด่นที่มีความต่อเนื่อง ส่วนทฤษฎีราฟเซตใช้ลดคุณสมบัติและสร้างกฎที่ใช้จำแนกการบุกรุก โดยได้ทำการทดสอบกับไฟล์ข้อมูลที่อยู่ในฐานข้อมูลของ KDD จากการทดลองพบว่าระบบตรวจจับการบุกรุกโดยใช้ราฟฟิชซีสามารถลดคุณสมบัติได้เป็นจำนวนมาก และสามารถสร้างกฎที่ใช้จำแนกการบุกรุกที่มีความแม่นยำสามารถอธิบายได้และมีเปอร์เซ็นต์ความแม่นยำและถูกต้องสูง งานที่ทำต่อไปคือทำการทดลองกับการบุกรุกระบบเครือข่ายแบบ Remote to User Attacks และ User to Root Attacks และดำเนินแก้ไข fuzzy space ให้เหมาะสมแต่ละคุณลักษณะเพื่อทำให้สามารถสร้างกฎที่ใช้จำแนกการบุกรุกที่สามารถตรวจจับการบุกรุกได้อย่างแม่นยำและหลากหลายกว่านี้

6. กิตติกรรมประกาศ

ขอขอบคุณ นายพรเทพ ไรจนวสุ นักศึกษาปริญญาเอก ภาควิชาวิศวกรรมคอมพิวเตอร์ สาขาวิชาวิศวกรรมไฟฟ้า คณะวิศวกรรมศาสตร์ สถาบันเทคโนโลยีพระจอมเกล้าพระนครเหนือ ที่มีส่วนช่วยเหลือในการให้คำแนะนำและตรวจทานบทความ

7. เอกสารอ้างอิง

- [1] M. Mohajerani, A. Moeini and M. Kianie, "NFIDS: A Neuro-fuzzy Intrusion Detection System", In Proceedings of the 10th IEEE International Conference on Electronics, Circuits and Systems, 2003, pp348-351.
- [2] W.Lee and S.J. Stolfo, "A Framework for Constructing Features and Models for Intrusion Detection Systems", In ACM Transactions on Information and System Security, Vol. 3, No. 4 ,November 2000, pp. 227-261
- [3] J. Gomez and D. Dasgupta, "Evolving Fuzzy Classifiers for Intrusion Detection", In Proceedings of the 2002 IEEE Workshop on Information Assurance United States Military Academy, West Point, NY June 2001
- [4] S. Lee, D. Heinbuch, "Training a Neural-Network Based Intrusion Detector to Recognize Novel attacks", Information Assurance and Security, pp.40-46,2000.
- [5] KDD Cup 1999: <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>
- [6] J.K. Liang, Y.Zhang, Y.B. Qu, "A Heuristic Algorithm of Attribute Reduction in Rough Set", In Proceedings of the Fourth International Conference on Machine Learning and Cybernetics, Guangzhou, 18-21 August 2005.
- [7] J. Huang, S. Li, C. Man, "A T-S Type of Rough Fuzzy Controller Based on Process Input-output Data", In Proceedings of the 42nd IEEE conference on Decision and control Maui, Hawaii USA, December 2003
- [8] Z. Lian-hua, Z. Guan-hua, Y. Lang, "Intrusion detection using rough set classification", In Journal of Zhejiang University SCIENCE, pp.1076-1086, 2004.

ประวัติผู้เขียน

ชื่อ-สกุล	นายธรรมกร ครองไตรภพ
วันเดือนปีเกิด	วันที่ 14 กรกฎาคม พ.ศ. 2512 ณ จังหวัดขอนแก่น
ที่อยู่	904 ถนน สุรนารายณ์ ตำบล ในเมือง อำเภอเมือง จังหวัด นครราชสีมา 30000
ประวัติการศึกษา	
พ.ศ. 2541	สำเร็จการศึกษาวิศวกรรมศาสตรบัณฑิต สาขาวิศวกรรมคอมพิวเตอร์ วิทยาเขตเทเวศร์ สถาบันเทคโนโลยีราชมงคล
พ.ศ. 2545	เข้าศึกษาต่อในระดับวิศวกรรมศาสตรมหาบัณฑิต สาขาวิศวกรรมคอมพิวเตอร์ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
ประสบการณ์ทำงาน	
พ.ศ. 2535-ปัจจุบัน	เข้ารับราชการตำแหน่งอาจารย์ประจำแผนกวิชาเทคนิคคอมพิวเตอร์ คณะวิชาไฟฟ้า สถาบันเทคโนโลยีราชมงคล วิทยาเขตภาคตะวันออกเฉียงเหนือ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้