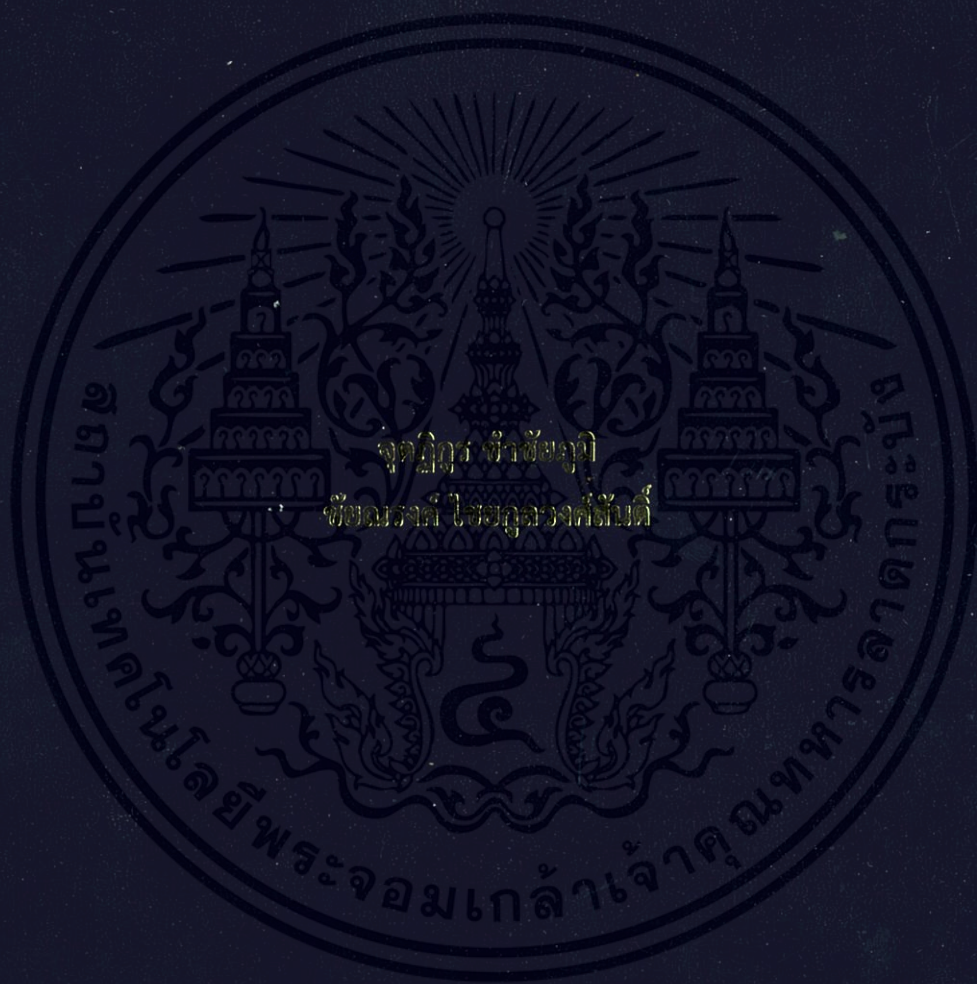


ระบบตรวจจับบอตเน็ตจากพฤติกรรมการใช้งานเครือข่าย  
NETWORK BEHAVIOR-BASED BOTNET DETECTION SYSTEM



ปริญญาโทนี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิศวกรรมศาสตรบัณฑิต

สาขาวิชาวิศวกรรมคอมพิวเตอร์

คณะวิศวกรรมศาสตร์

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

ปีการศึกษา 2556

ระบบตรวจจับบอทเน็ตจากพฤติกรรมการใช้งานเครือข่าย  
NETWORK BEHAVIOR-BASED BOTNET DETECTION SYSTEM



จตุฎฎิฎฎุฎ ฎฎฎฎฎฎ  
ฎฎฎฎฎฎ ฎฎฎฎฎฎ

ปริญญานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิศวกรรมศาสตรบัณฑิต

สาขาวิชาวิศวกรรมคอมพิวเตอร์

คณะวิศวกรรมศาสตร์

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษานานาชาติ ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งยังขอสงวนลิขสิทธิ์ในสิ่งพิมพ์ฉบับนี้ของมหาวิทยาลัยเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

ปีการศึกษา 2556

ปริญญาานิพนธ์ปีการศึกษา 2556

สาขาวิชาวิศวกรรมคอมพิวเตอร์

คณะวิศวกรรมศาสตร์ สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง

เรื่อง ระบบตรวจจับบอทเน็ตจากพฤติกรรมการใช้งานเครือข่าย

NETWORK BEHAVIOR-BASED BOTNET DETECTION SYSTEM

ผู้จัดทำ

1. นายจตุภูมิ ขำชัยภูมิ รหัสนักศึกษา 53010247

2. นายชัยณรงค์ ไชยกุลวงศ์สันติ รหัสนักศึกษา 53010335



.....อาจารย์ที่ปรึกษา

(ดร. ชาญชัย ตรีภาค)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

# ระบบตรวจจับบอทเน็ตจากพฤติกรรมการใช้งานเครือข่าย

นายจตุภูมิ ขำชัยภูมิ 53010247  
นายชัยณรงค์ ไชยกุลวงศ์สันต์ 53010335  
ดร. ธนัญชัย ตรีภาค อาจารย์ที่ปรึกษา  
ปีการศึกษา 2556

## บทคัดย่อ

ในปัจจุบันเครื่องคอมพิวเตอร์ของผู้ใช้ตามบ้านหรือในองค์กรกำลังถูกแฮคเกอร์เข้าควบคุมและแปรสภาพเครื่องดังกล่าวให้ทำงานเป็น “bots”, “zombie” หรือ “drones” กล่าวคือทำการติดตั้งซอฟต์แวร์ที่ทำให้แฮคเกอร์สามารถควบคุมระยะไกลได้ โดยจะเรียกเครื่องคอมพิวเตอร์ต่างๆ ที่ถูกแฮคเกอร์ควบคุมว่า “บอทเน็ต (Botnet)” ซึ่งวัตถุประสงค์ในการรวบรวมบอทเน็ตคือเพื่อใช้ในการประกอบกิจกรรมที่ขัดต่อกฎหมาย อย่างเช่น ส่งสแปมเมลล์ หรือ โจมตีโดยวิธี Denial of Service (DOS attack) จากปัญหาดังกล่าวผู้พัฒนาจึงทำการศึกษาและทำความเข้าใจกับบอทเน็ตรวมถึงรูปแบบพฤติกรรมที่บอทเน็ตกระทำกับเครื่องที่ตกเป็นเป้าหมาย ทั้งระบบการทำงานต่างๆ ของบอทเน็ต และพัฒนาโปรแกรมสำหรับวิเคราะห์และตรวจหาทราบฟิคที่อาจจะเป็นของบอทเน็ตในเครือข่ายที่ดูแลได้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

# NETWORK BEHAVIOR-BASED BOTNET DETECTION SYSTEM

Mr. Juttikhun Khamchaiyaphum 53010247

Mr. Chainarong Chaikoonvongsun 53010335

Dr. Thanunchai Threepak Advisor

Academic Year 2013

## ABSTRACT

Home or office computers are under threats from hackers who take over the control and alter the system to operate as "bots", "zombie", or "drones" by installing some remote control software enabling the hackers to control from afar. The controlled computers will be called "Botnet" and used by hackers collectively for illegal purposes such as sending spam mails or attacking by denial of service. From such problems, our team of developers have studied and understood about Botnet, including behavioral pattern that Botnet has on the victimized system, various working system of Botnet, and then has developed an application for analysis and searching for traffic that has a possibility to be Botnet in the network under supervision.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



# สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	I
บทคัดย่อภาษาอังกฤษ.....	II
กิตติกรรมประกาศ.....	III
สารบัญ.....	IV
สารบัญตาราง.....	VII
สารบัญรูป.....	VIII
บทที่ 1 บทนำ.....	1
1.1 ความสำคัญและที่มาของโครงการ.....	1
1.2 วัตถุประสงค์ของโครงการ.....	1
1.3 ขอบเขตของโครงการ.....	2
1.4 วิธีการดำเนินการ.....	2
1.5 ประโยชน์ที่คาดว่าจะได้รับ.....	2
1.6 ส่วนประกอบของปริญญานิพนธ์.....	2
บทที่ 2 ทฤษฎีที่เกี่ยวข้อง.....	4
2.1 บอทเน็ต.....	4
2.2 ลักษณะของบอทเน็ต.....	4
2.2.1 การแพร่กระจายตนเอง.....	5
2.2.2 พฤติกรรมการโจมตี.....	6
2.2.3 โครงสร้างในการใช้คำสั่งสื่อสารและควบคุมบอท.....	7
2.2.4 โพรโทคอลในการสื่อสาร.....	10
2.2.5 กลไกการควบคุมและสั่งการศูนย์กลาง.....	11
2.3 เหตุผลที่เลือกเอชทีทีพีบอทเน็ตในการวิจัย.....	12
2.4 วิธีการที่ใช้ในการตรวจสอบบอทเน็ต.....	13
2.4.1 อันนี้พ็อตและอันนี้เน็ต.....	13
2.4.2 ตรวจจับโดยลายเซ็น.....	13

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดลอกแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## สารบัญ (ต่อ)

	หน้า
2.4.3 ตรวจสอบโดยการตรวจสอบดีเอ็นเอส.....	13
2.5 การตรวจสอบโดยการวิเคราะห์พฤติกรรมของเครือข่าย.....	14
2.5.1 ตัวอย่างการตรวจสอบโดยการวิเคราะห์พฤติกรรมตามช่วงของเวลา (Period).....	14
2.5.2 เหตุผลที่เลือกวิธีการตรวจสอบโดยใช้การวิเคราะห์พฤติกรรมเครือข่าย.....	15
2.6 เทคนิคการแบ่งกลุ่มแบบเชิงลำดับชั้น.....	16
2.6.1 ประโยชน์ของอัลกอริทึมเบิร์ช.....	16
2.6.2 หลักการของอัลกอริทึมเบิร์ช.....	16
2.7 ระยะห่างเชิงเลขคณิต, ระยะห่างเชิงประเภท และระยะห่างเชิงลำดับชั้นสำหรับ อัลกอริทึมเบิร์ช.....	20
2.7.1 คุณสมบัติที่จำเป็นสำหรับการวัดความคล้ายคลึงระหว่างข้อมูลแบบลำดับชั้น.....	23
2.7.2 การวิเคราะห์หมายเลขไอพีแอดเดรสแบบสมบัติเชิงลำดับชั้น.....	25
2.7.3 การรวมผลต่างระยะห่างเชิงเลขคณิต, เชิงประเภท และเชิงลำดับชั้น.....	26
2.7.4 การคำนวณรัศมีของกลุ่มข้อมูล.....	26
2.8 เน็ตโพลว์.....	27
2.9 เครื่องมือที่ใช้ในการดำเนินงาน.....	28
บทที่ 3 การออกแบบและการพัฒนา.....	29
3.1 รายละเอียดของระบบที่พัฒนา.....	29
3.1.1 รายละเอียดการนำเข้าข้อมูล (Input Specification).....	29
3.1.2 รายละเอียดผลลัพธ์ของระบบ (Output Specification).....	29
3.1.3 ขอบเขตของระบบที่พัฒนา.....	30
3.1.4 ข้อจำกัดของระบบที่พัฒนา.....	30
3.1.5 เครื่องมือที่ใช้ในการพัฒนา.....	30
3.2 รูปแบบโครงสร้างของระบบ.....	30
3.2.1 ส่วนรวบรวมข้อมูล.....	31
3.2.2 ส่วนเก็บรักษาข้อมูล.....	32

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งยังมีให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## สารบัญ (ต่อ)

	หน้า
3.2.3 ส่วนวิเคราะห์ข้อมูล .....	33
3.3 แนวคิดการวิเคราะห์ข้อมูล .....	34
3.3.1 รูปแบบการแบ่งกลุ่มข้อมูล .....	35
3.3.2 การสรุปผล .....	36
3.4 ส่วนของผู้ใช้งานและยูสเคสไดอะแกรม .....	37
บทที่ 4 การทดลองและผลการทดลอง .....	42
4.1 การทดลองโปรแกรมเอฟโพรบ (Fprobe) .....	42
4.1.1 การติดตั้งและตั้งค่าเอฟโพรบ (Fprobe) และการเก็บข้อมูล NetFlow .....	42
4.2 การทดลองการทำงานของอัลกอริทึมเบิร์ช .....	44
4.3 การทดลองการทำงานของ Botnet (VertexNet รุ่น 1.2.1) .....	44
4.3.1 ทดลองการติดตั้งบอท .....	44
4.3.2 ทดลองการสั่งงานบอทและการสังเกตผล .....	45
4.4 การทดลองการทำงานของ Botnet (Zeus รุ่น 2.1.0.1) .....	48
4.4 การทดลองการวิเคราะห์การสื่อสารของบอทกับเซิร์ฟเวอร์ .....	49
4.4.1 พฤติกรรมการเชื่อมต่อเป็นคาบของบอทเน็ต Zeus รุ่น 2.1.0.1 และ VertexNet รุ่น 1.2.1 .....	49
4.5 การทดลองการคำนวณตามเอกสารอ้างอิง [1] .....	51
บทที่ 5 บทสรุปและข้อเสนอแนะ .....	53
5.1 บทสรุป .....	53
5.2 ปัญหา อุปสรรค และแนวทางการแก้ไข .....	53
5.3 แนวทางในการพัฒนาต่อ .....	54

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
บรรณานุกรม .....

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

# สารบัญตาราง

ตารางที่

หน้า

2.1 ตารางเวลาในแต่ละช่วงเวลาที่มีกิจกรรมติดต่อกับซีแอนซีเซิร์ฟเวอร์ ..... 15



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

# สารบัญรูป

รูปที่	หน้า
2.1 ตัวอย่างพีชชิงอีเมลล์.....	5
2.2 ตัวอย่างโครงสร้างสตาร์ของซีแอนซีบอทเน็ต .....	8
2.3 ตัวอย่างโครงสร้างแบบมัลติเซิร์ฟเวอร์ของซีแอนซีบอทเน็ต .....	9
2.4 ตัวอย่างโครงสร้างแบบลำดับชั้นของซีแอนซีบอทเน็ต .....	9
2.5 ตัวอย่างโครงสร้างแบบสุ่มของซีแอนซีบอทเน็ต.....	10
2.6 ไออาร์ซีเบสซีแอนซีบอทเน็ต.....	11
2.7 เอชทีทีพีเบสซีแอนซีบอทเน็ต .....	12
2.8 โค้ดแสดงกลไกการใช้คำสั่ง sleep() อย่างง่าย .....	14
2.9 การแบ่งกลุ่มแบบเชิงลำดับชั้น .....	16
2.10 ซีเอฟทีรีของแต่ละโหนดย่อย .....	17
2.11 ตัวอย่างซีเอฟทีรี.....	18
2.12 ตัวอย่างอื่นๆของซีเอฟทีรี .....	19
2.13 ตัวอย่างต้นไม้ทวิภาคสำหรับหมายเลขพอร์ต.....	23
2.14 ตัวอย่างข้อมูลที่จัดเรียงเชิงลำดับชั้น.....	24
3.1 ภาพรวมโครงสร้างของระบบ .....	31
3.2 ผังแสดงการทำงานของส่วนรวบรวมข้อมูล .....	32
3.3 ผังแสดงการทำงานของส่วนเก็บรักษาข้อมูล .....	33
3.4 ผังแสดงการทำงานของส่วนวิเคราะห์ข้อมูล .....	34
3.5 โค้ดเทียมสำหรับอัลกอริทึมเพิ่มเรคคอร์ด X ในซีเอฟทีรี .....	36
3.6 ยูสเคสไดอะแกรม .....	38
3.7 หน้าหลักของระบบซึ่งเป็นหน้าสำหรับการจัดการอุปกรณ์.....	39
3.8 หน้าเพิ่มไอพีแอดเดรสของเราท์เตอร์ให้กับโปรแกรม .....	39
3.9 หน้าสำหรับตั้งค่าอุปกรณ์.....	40
3.10 หน้าสำหรับตั้งค่าอุปกรณ์ขั้นสูง .....	40
3.11 หน้าการตั้งค่าการแสดงผล .....	41
3.12 หน้ารายงานจำนวนการเชื่อมต่อและกราฟแสดงผลจำนวนบอทเน็ตที่มีในเครือข่าย .....	41

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งยังมีให้อัปโหลดเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## สารบัญรูป (ต่อ)

รูปที่	หน้า
4.1 แสดงการรันคำสั่ง <code>fprobe -l eth0 ip -d 15 -e 60 161.246.5.14:5555</code> .....	42
4.2 รูปภาพแสดงการทดลองเก็บข้อมูล NetFlow.....	43
4.3 รูปภาพแสดงฐานข้อมูล NetFlow.....	43
4.4 ผลลัพธ์การแบ่งกลุ่มตามอัลกอริทึม .....	44
4.5 การเชื่อมต่อระหว่าง C&C กับบอทโดยที่จะมีค่า “idle” ซึ่งเป็นค่าที่บอกว่าบอทนั้นมีการเคลื่อนไหวหรือไม่ ถ้าหากมีการเคลื่อนไหวของบอทเช่นผู้ใช้มีการขยับเมาส์บอทก็จะส่งแพ็คเก็ตเพื่อรีเซ็ตค่า “idle” เป็น 0 .....	45
4.6 การทดลองสั่งงานบอทด้วยคำสั่ง <code>getproc</code> .....	46
4.7 ผลลัพธ์ของคำสั่ง <code>getproc</code> ซึ่งจะเห็นได้ว่าการรีเทิร์นโพรเซสที่ทำงานอยู่บนเครื่องที่เป็นบอท .....	47
4.8 การทดลองการดักจับแพ็คเก็ตเพื่อสังเกตการณ์การทำงานของ C&C โดยเป็นการดักจับแพ็คเก็ตขณะที่มีการสั่งคำสั่ง <code>shutdown -r -t 0</code> .....	47
4.9 ผลที่ได้จากการดักจับแพ็คเก็ตของ Wireshark ที่เป็นของบอท Zeus ที่ร้องขอไปที่ <code>/cfg.bin</code> ....	48
4.10 หน้าเว็บของบอทมาสเตอร์ที่คอยควบคุมและสั่งคำสั่งให้กับบอทลูก .....	49
4.11 รูปแสดงการกรองแพ็คเก็ตจากหมายเลขไอพี, หมายเลขพอร์ต, โพรโตคอลและปริมาณข้อมูล .....	50
4.12 ผลที่ได้จากการกรองแพ็คเก็ตที่มีการส่งคำร้องขอไปที่ VertexNet เซิร์ฟเวอร์โดยมีการเชื่อมต่อส่งไปทุกๆ 30 วินาที.....	50
4.13 ผลการคำนวณ CF Entry ของข้อมูลแต่ละชุด .....	51
4.14 การจำลอง CF Tree ที่ได้จากการคำนวณ.....	52

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

# บทที่ 1

## บทนำ

### 1.1 ความสำคัญและที่มาของโครงการ

ปัจจุบันเครือข่ายคอมพิวเตอร์และอินเทอร์เน็ตมีการใช้งานเพิ่มขึ้นอย่างรวดเร็ว ระบบรักษาความปลอดภัยจึงเป็นเรื่องสำคัญอย่างมาก เพราะข้อมูลต่างๆ อาจเกิดการสูญหายหรือมีผู้อื่นทำการขโมยข้อมูลเหล่านั้นไป โดยใช้การโจมตีต่างๆ ซึ่งการโจมตีหรือภัยคุกคามมีรูปแบบที่แตกต่างทำให้สามารถจำแนกรูปแบบการโจมตีหรือภัยคุกคามเหล่านี้ได้ คือ Virus, Worm, Trojans, Backdoor, Spyware และ Phishing Botnets ซึ่งจากรูปแบบการโจมตีดังกล่าว บอทเน็ตนั้นถือว่าการโจมตีที่อันตรายและสร้างความเสียหายมากที่สุด

บอทเน็ตเป็นการโจมตีที่มีลักษณะวิธีการแพร่กระจายตนเองที่หลากหลาย ส่งผลให้เครื่องคอมพิวเตอร์ที่ถูกโจมตีทำงานภายใต้อำนาจควบคุมของผู้อื่น เมื่อได้รับคำสั่งจากผู้ควบคุม หรือเรียกว่า บอทมาสเตอร์ เพื่อให้คอมพิวเตอร์ทำงานในลักษณะต่างๆ ให้บรรลุเป้าหมายตามที่ต้องการ โดยมีกลไกในการควบคุมและสั่งการที่แตกต่างกัน (เช่น IRC, HTTP, P2P) จะพบว่าบอทเน็ตไม่ได้เป็นเพียงภัยคุกคามต่อระบบเครือข่ายและอินเทอร์เน็ตเท่านั้น แต่ยังมีส่วนร่วมในรูปแบบอื่นๆ ของภัยคุกคามและการโจมตีอีกด้วย จากปัญหาต่างๆ ที่เกิดจากภัยคุกคามของบอทเน็ต ทำให้จุดประสงค์ในการพัฒนาโครงการนี้คือ มีระบบที่สามารถตรวจจับบอทเน็ตโดยใช้หลักการวิเคราะห์พฤติกรรมของบอท ซึ่งในขั้นตอนการดำเนินงาน จะรวบรวมข้อมูลต่างๆ ในเครือข่าย แล้วทำการเก็บรวบรวมไว้ในฐานข้อมูลของเซิร์ฟเวอร์ จากนั้นจะนำข้อมูลมาตรวจสอบโดยใช้หลักการแบ่งกลุ่มเชิงลำดับชั้น (Hierarchical Clustering) เพื่อแบ่งลักษณะการใช้งานระหว่างกลุ่มพฤติกรรมของบอทกับกลุ่มพฤติกรรมของบุคคลใช้งานทั่วไปออกจากกัน

### 1.2 วัตถุประสงค์ของโครงการ

- 1) เพื่อพัฒนาระบบในการตรวจจับบอทต่างๆ โดยใช้หลักการวิเคราะห์พฤติกรรมของบอท
- 2) เพื่อให้ผู้ใช้ระบบสามารถล่วงรู้ทันเหตุการณ์ว่าเครื่องคอมพิวเตอร์นี้ถูกคุกคามโดยบอท

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

### 1.3 ขอบเขตของโครงการ

- 1) พัฒนาแอปพลิเคชันสำหรับวิเคราะห์กราฟฟิคที่ผิดปกติ โดยใช้สถิติของพฤติกรรมการใช้งานเครือข่าย เพื่อตรวจจับบอทเน็ตที่ทำงานบนวินโดวส์เซิร์ฟเวอร์
- 2) ทำการตรวจจับบอทเน็ตที่ทำงานบนระบบเครือข่าย

### 1.4 วิธีการดำเนินการ

- 1) ศึกษาการทำงานของบอทเน็ต
- 2) ศึกษารูปแบบการตรวจจับแบบวิเคราะห์พฤติกรรม
- 3) ศึกษางานวิจัยของ Meisam Eslahi เรื่อง Improving HTTP-Based Botnet Detection by Using Network Behavior Analysis System
- 4) ศึกษาวิธีการจัดแบ่งกลุ่มของบอท Hierarchical Clustering (An Efficient Clustering Scheme to Exploit Hierarchical Data in Network Traffic Analysis : BIRCH)
- 5) จำลองการทำงานของบอทและจัดเก็บสถิติไว้ใช้ในการจัดกลุ่ม
- 6) พัฒนาโปรแกรมสำหรับคำนวณการจัดแบ่งกลุ่มของบอท
- 7) ทดลองใช้งานโปรแกรมที่พัฒนาขึ้นมา
- 8) สรุปผลการทดลองและแก้ไขปรับปรุงกรณีมีข้อผิดพลาด

### 1.5 ประโยชน์ที่คาดว่าจะได้รับ

- 1) ได้รับความรู้เกี่ยวกับบอทแต่ละประเภท
- 2) ได้รับความรู้เกี่ยวกับกระบวนการจัดแบ่งกลุ่มของบอทชนิดต่างๆ
- 3) ได้นำโปรแกรมตรวจจับบอทเน็ตไปใช้งานจริงในการตรวจจับบอทเน็ตบนระบบต่างๆ

### 1.6 ส่วนประกอบของปฏิญานิพนธ์

ปฏิญานิพนธ์ฉบับนี้ได้แบ่งเนื้อหาออกเป็น 5 บท โดยมีรายละเอียดดังต่อไปนี้

บทที่ 1 บทนำ กล่าวถึง ความสำคัญและที่มาของโครงการ วัตถุประสงค์ของโครงการ

ขอบเขตของโครงการ วิธีการดำเนินการ ประโยชน์ที่คาดว่าจะได้รับ และส่วนประกอบของปฏิญานิพนธ์

บทที่ 2 ทฤษฎีที่เกี่ยวข้อง กล่าวถึง ทฤษฎีพื้นฐานที่ใช้ในโครงการ ประกอบด้วย การแบ่งกลุ่มแบบเชิงลำดับขั้นด้วยกระบวนการของอัลกอริทึมเบิร์ชและการวิเคราะห์พฤติกรรมของเครือข่าย นำไปใช้

บทที่ 3 การออกแบบและการพัฒนา กล่าวถึง รายละเอียดระบบที่พัฒนารูปแบบโครงสร้างของระบบและแนวคิดการวิเคราะห์ข้อมูล

บทที่ 4 การทดลองและผลการทดลองกล่าวถึง การทดลองโปรแกรมและการทดลองการทำงานของอัลกอริทึมที่ใช้

บทที่ 5 บทสรุป กล่าวถึงบทสรุปของโครงการ ปัญหาและอุปสรรคต่างๆของโครงการ แนวทางการแก้ไขและแนวทางการพัฒนาต่อไป



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## บทที่ 2

# ทฤษฎีที่เกี่ยวข้อง

### 2.1 บอทเน็ต

บอทเป็นโปรแกรมอัตโนมัติที่สามารถดำเนินการและทำงานด้วยคำสั่งซ้ำๆ กันได้อย่างรวดเร็ว มีจุดประสงค์หลักเพื่อดำเนินกิจกรรมที่เป็นอันตรายในเครือข่ายคอมพิวเตอร์ บอทจำนวนมากจะแพร่กระจายตนเองไปยังเครื่องคอมพิวเตอร์ต่างๆ และเชื่อมต่อกับผู้ใช้อื่นๆ ผ่านทางการใช้งานอินเทอร์เน็ต เมื่อมีบอทจำนวนมากทำงานในลักษณะเป็นกลุ่ม จะเรียกพฤติกรรมดังกล่าวว่า “บอทเน็ต (Botnet)” เครือข่ายของบอทหรือบอทเน็ตจะมีขนาดเครือข่ายตั้งแต่ขนาดเล็กๆ ระดับเพียงพันบอท จนถึงเครือข่ายขนาดใหญ่ซึ่งอาจมีจำนวนบอทมากถึงล้านบอท

บอทถูกออกแบบมาเพื่อทำให้เครื่องคอมพิวเตอร์ในเครือข่ายติดไวรัสหรือถูกคุกคาม เครื่องคอมพิวเตอร์ที่ถูกคุกคามจะกลายเป็นส่วนหนึ่งของบอทโดยที่เจ้าของนั้นไม่รับรู้ถึงภัยคุกคามที่เกิดขึ้นได้เลย ผู้ที่ควบคุมบอทเหล่านี้ เรียกว่า บอทมาสเตอร์ (Bot Master) จะทำการส่งคำสั่งประสงค์ร้ายต่างๆ ไปยังบอททั้งหมดผ่านทางอินเทอร์เน็ตและเซิร์ฟเวอร์ ที่มีชื่อว่า command and control (เรียกย่อๆว่า ซีแอนด์ซี) เซิร์ฟเวอร์

วงจรชีวิตของบอทเน็ตสามารถแบ่งออกได้เป็น 4 ช่วง โดยช่วงสัปดาห์แรกของการถูกคุกคาม เป้าหมายหลักของบอท คือ การรวบรวมข้อมูลที่สามารถเข้าถึงได้ทั้งหมดของเครื่องที่ถูกคุกคาม เช่น รหัสผ่านทางเว็บไซต์, อีเมลล์, บัญชีธนาคาร เป็นต้น โดยจะใช้วิธีการต่างๆ ในการรวบรวมข้อมูล เช่น การเปลี่ยนเส้นทางเดินของข้อมูลไปยังฟิชซิงเว็บไซต์ (phishing websites) เป็นต้น ในช่วงสัปดาห์ที่ 2 บอทจะรวบรวมข้อมูลในส่วนของโลคอลเน็ตเวิร์ค (Local Network) เช่น ข้อมูลการใช้งานเครือข่ายในองค์กร หรือบริษัท ช่วงสัปดาห์ที่ 3 จะเป็นช่วงแพร่กระจายตนเองของบอท เช่น การส่งสแปม หรือ การแนบตนเองไปพร้อมกับการส่งข้อความผ่านทางสื่อออนไลน์ต่างๆ และช่วงสัปดาห์ที่ 4 บอททั้งหมดจะถูกควบคุมโดยบอทตัวหลัก เพื่อรอรับคำสั่งสำหรับการโจมตีต่อไป

### 2.2 ลักษณะของบอทเน็ต

ลักษณะหลักของบอทเน็ต ระบุโดย เทรินไมโคร (Trend Micro) และ Jose Nazario มีดังต่อไปนี้

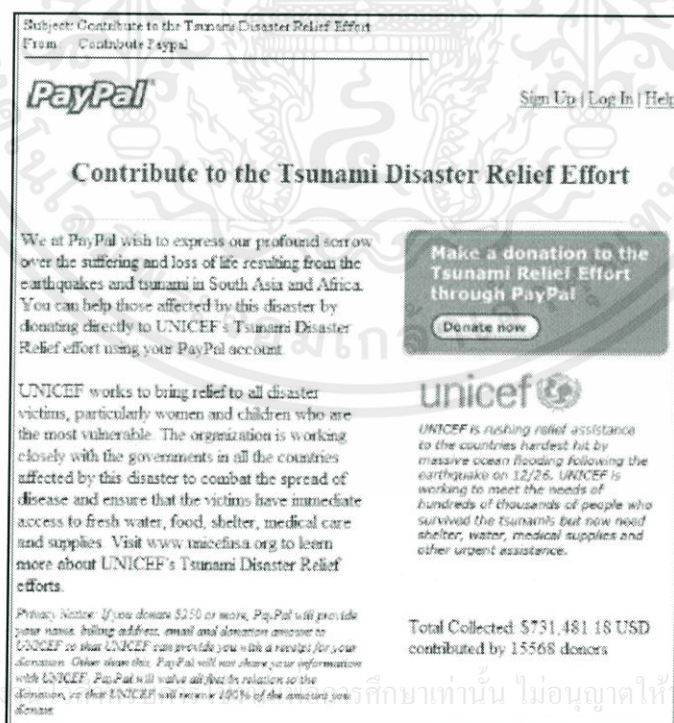
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## 2.2.1 การแพร่กระจายตนเอง

บอทเน็ตมีกรรมวิธีที่หลากหลายในการแพร่กระจายตนเองไปยังเครื่องคอมพิวเตอร์ โดยทั่วไปเรียกว่า เวิร์ม (Worm) ซึ่งลักษณะการแพร่กระจายตนเองของบอทมักจะมีพฤติกรรมแบบซ้ำ เดิม Jose Nazario จึงแบ่งประเภทของการแพร่กระจายออกเป็น 4 รูปแบบดังต่อไปนี้

### 1) อีเมลล์ (Email)

อีเมลล์จำนวนมากอาจมีไฟล์บางไฟล์ที่บอทแฝงตัวแนบไปกับอีเมลล์ เพื่อจะทำการแพร่กระจายตนเอง หากไฟล์นั้นถูกเปิดใช้งาน อาจทำให้เครื่องคอมพิวเตอร์เครื่องนั้นถูกควบคุมได้ทันที ในอดีตนั้นเราสามารถทำการป้องกันได้ง่าย เนื่องจากอีเมลล์จำพวกนี้มักถูกส่งจากบุคคลที่เราไม่รู้จัก แต่ในปัจจุบันอีเมลล์เหล่านี้ใช้ความรู้ทางวิศวกรรมสังคม (Social Engineering) แฝงตัวเข้ามาโดยอาจอ้างว่าอีเมลล์นี้ถูกส่งมาจากแหล่งที่มาที่มีความน่าเชื่อถือ เช่น อีเมลล์มาจากธนาคาร หรือ บุคคลที่รู้จัก บ่อยครั้งที่เครื่องคอมพิวเตอร์ที่ถูกบุคคลอื่นขโมยข้อมูลที่อยู่ และส่งต่ออีเมลล์ไปยังที่อยู่อื่นๆ เพื่อทำการแพร่กระจายบอทออกไป รูปที่ 2.1 แสดงตัวอย่าง พิซซิ่งอีเมลล์ ที่ส่งไปยังเครื่องเหยื่อ โดยอ้างว่าเป็นการระดมเงินสำหรับภัยพิบัติสึนามิ จะเห็นได้ว่าเป็นเรื่องที่ยากที่จะสามารถแยกความแตกต่างระหว่างอีเมลล์จริงออกจากพิซซิ่งอีเมลล์จำพวกนี้



รูปที่ 2.1 ตัวอย่างพิซซิ่งอีเมลล์

## 2) ลิงค์สแปม (Link Spam)

ลิงค์สแปม คือ ยูอาร์แอล (URL) ที่มีเนื้อหาที่เป็นอันตรายซึ่งอาจได้มาโดยหลากหลายกรรมวิธี เช่น อีเมลล์, เว็บไซต์ออนไลน์ต่างๆ และแชทต่างๆ ทั้งนี้ทุกกรรมวิธีต้องมีการพยายามที่จะโน้มน้าวให้ผู้ใช้นั้นคิดว่าเป็นความจริง ข้อความเหล่านี้มักจะปรากฏขึ้นมาจากเพื่อนหรือองค์กรที่เชื่อถือได้ แต่ความจริงแล้วได้ถูกปลอมแปลงมาจากอีกบุคคลหนึ่ง

## 3) เว็บไซต์ (Website)

เว็บไซต์ที่ปลอมแปลงนั้นจะมีการติดตั้งตามเว็บที่เป็นที่รู้จักดี เช่น ยูทูบ หรือ เว็บไซต์ที่เชื่อถือได้ที่มีการถูกแฮคและถูกฝังโค้ดที่เป็นอันตราย สามารถแบ่งเป็นการโจมตีได้ 2 ประเภท คือ การโจมตีฝั่งผู้ใช้ และการใช้ประโยชน์จากการดาวน์โหลด ในส่วนของการโจมตีฝั่งผู้ใช้นั้น ผู้ใช้จะมีการใช้บริการเว็บไซต์ต่างๆ ที่ถูกฝังโค้ดที่เป็นอันตราย โค้ดจะทำงานโดยพยายามที่จะใช้ช่องโหว่ต่างๆ ของเว็บเบราว์เซอร์ในการเข้าถึงเครื่องของผู้ใช้งาน หากประสบความสำเร็จ เครื่องของผู้ใช้งานจะถูกควบคุมโดยผู้อื่นทันที ส่วนการใช้ประโยชน์จากการดาวน์โหลดนั้น เมื่อเกิดการแจ้งเตือนในการดาวน์โหลดไฟล์ แต่ผู้ใช้นั้นยอมรับการดาวน์โหลดนั้น คอมพิวเตอร์ทำการรันไฟล์และส่งผลให้เครื่องคอมพิวเตอร์ถูกคุกคามในทันที

## 4) Exploits

วิธีการนี้จะคล้ายคลึงกับการโจมตีฝั่งผู้ใช้ แต่การใช้ช่องโหว่นั้นไม่ได้จำกัดว่าจะต้องมาจากเว็บเบราว์เซอร์เท่านั้น เมื่อเครื่องคอมพิวเตอร์ถูกคุกคาม ก็จะพยายามสแกนเครือข่ายเพื่อหาเครื่องคอมพิวเตอร์เครื่องอื่นๆ และจะทำการโจมตีหรือหารหัสผ่านของเครื่องที่ง่ายต่อการเข้าถึงต่อไป

### 2.2.2 พฤติกรรมการโจมตี

บอทเน็ตจะถูกสร้างขึ้นโดยผู้สร้างซึ่งมีการนำไปใช้งานในลักษณะที่แตกต่างกันดังต่อไปนี้

#### 1) สแปมและฟิชชิงอีเมลล์

ในปี 2009 มีอีเมลล์ที่เป็นสแปมมากถึงแสนล้านฉบับถูกส่งทุกๆ วัน คิดเป็น 87.7 % ของอีเมลล์ทั้งหมด ซึ่ง 83% ของอีเมลล์สแปมนั้นถูกส่งมาจากบอทเน็ต สแปมเป็นปัญหาใหญ่ที่ไม่ใช่แค่เพียงสร้างความรำคาญให้กับผู้ใช้งานเพียงอย่างเดียว แต่ยังสามารถทำฟิชชิงข้อมูลต่างๆ ที่เป็นความลับได้อีกด้วย เช่น ชื่อผู้ใช้, รหัสผ่าน, หมายเลขบัตรเครดิต เป็นต้น

#### 2) ดีดีโอเอส (DDoS)

การโจมตีแบบดีดีโอเอส (Distributed Denial of Service) เป็นปัญหาสำคัญอีกประการหนึ่งเพราะยากต่อการป้องกัน การโจมตีแบบดีดีโอเอส (DoS) จะเกิดขึ้นเมื่อเซิร์ฟเวอร์

ถูกร้องขอให้สร้างการเชื่อมต่อเป็นจำนวนมากจากโคลเอนท์ ทำให้มีปริมาณการไหลเวียนของข้อมูลบนเครือข่ายจำนวนมากจนกระทั่งเครื่องเซิร์ฟเวอร์ไม่สามารถตอบสนองได้ทั้งหมด ซึ่งแต่ก่อนนั้นการโจมตีชนิดนี้จะง่ายต่อการป้องกัน เพียงแค่ไฟร์วอลล์ในการบล็อกไอพีแอดเดรสของผู้โจมตี แต่ในกรณีที่มีการร้องขอการเชื่อมต่อมาจากโคลเอนท์หลายทีก็เป็น การยากที่จะทำการป้องกันได้อย่างเต็มประสิทธิภาพ

### 3) คีย์ล็อกกิง (Key Logging)

บอทเน็ตจะได้รับข้อมูลจำนวนมากจากคอมพิวเตอร์ที่ถูกคุกคาม ข้อมูลทุกส่วนในเครือข่ายหรือแม้แต่ขณะที่พิมพ์ข้อมูลบางอย่างบนเครื่องคอมพิวเตอร์ก็สามารถถูกคุกคามได้ทันที หากเครื่องคอมพิวเตอร์นั้นมีการใช้งานคีย์ล็อกกิง ซึ่งเป็นซอฟต์แวร์ที่จะบันทึกการกดแป้นพิมพ์ทุกอย่าง ไม่ว่าจะเป็นรหัสผ่านของเว็บไซต์ เช่น PayPal หรือ ธนาคาร เป็นต้น

### 4) Click Fraud

Click Fraud เป็นลักษณะการโกงของบอท เมื่อบอทเข้าไปแฝงตัวอยู่ในเว็บไซต์ต่างๆ จะสร้างทำเป็นผู้ใช้ทั่วไป และจะคลิกโฆษณาต่างๆ โดยอ้างว่าทำกำไรให้กับเจ้าของเว็บไซต์ ซึ่งการคลิกนี้จะไปนำเงินจากบริษัทโฆษณาออนไลน์ที่จ่ายเงินให้กับการคลิกโฆษณาแต่ละครั้ง ซึ่งเป็นการยากที่จะตรวจจับ เพราะทุกครั้งที่มีการคลิกนั้นจะมาจากต้นทางที่แตกต่างกัน เพราะมีการกระจายตัวไปทั่วอินเทอร์เน็ต

### 5) มัลแวร์ (Malware)

บอทเน็ตจะแพร่กระจายมัลแวร์เพื่อเพิ่มจำนวนเครื่องที่ถูกคุกคาม ซึ่งวิธีที่นิยมใช้คือผ่านทางอีเมลล์ แต่ในปัจจุบันนี้โซเชียลเน็ตเวิร์คและแอปพลิเคชันในการแชทต่างๆ เช่น เฟสบุ๊ค, ทวิตเตอร์ อาจมีการส่งข้อความที่ประกอบด้วยยูอาร์แอล ซึ่งจะนำไปสู่เว็บไซต์ที่อันตรายต่อผู้ใช้งาน

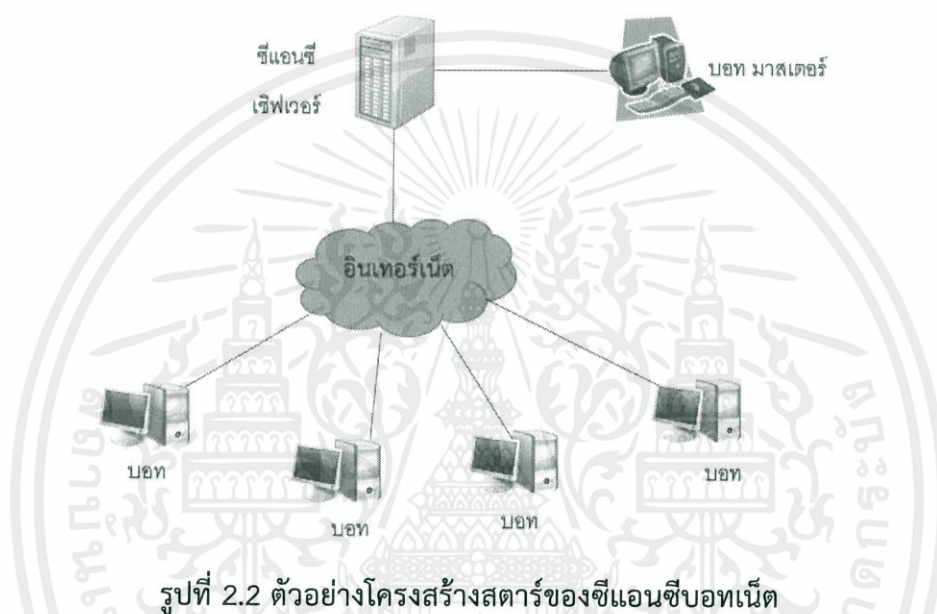
## 2.2.3 โครงสร้างในการใช้คำสั่งสื่อสารและควบคุมบอท

ซีแอนด์ซี (Command and Control) จะคอยสั่งการและควบคุมการสื่อสารระหว่างบอทกับบอทมาสเตอร์ บอททั้งหมดในเครือข่ายจะติดต่อกันผ่านทางมาสเตอร์เซิร์ฟเวอร์ โดยการส่งข้อมูลที่เก็บรวบรวมได้ (collected information) เช่น รหัสผ่าน มายิงมาสเตอร์เซิร์ฟเวอร์ นอกจากนี้ยังสามารถรับคำสั่งอื่นๆ จากเซิร์ฟเวอร์ได้ด้วย โครงสร้างของซีแอนด์ซีที่ถูกใช้โดยบอทเน็ตสามารถแบ่งได้เป็น 4 ประเภทดังนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆ ทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

### 1) สตาร์ (Star)

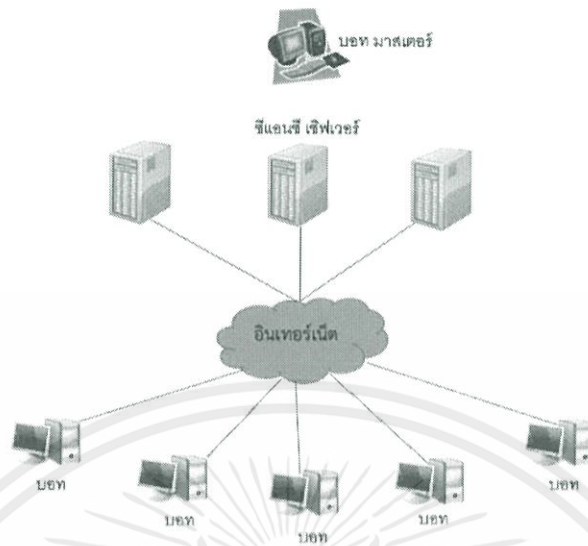
รูปแบบสตาร์จะทำให้บอททั้งหมดสื่อสารกันผ่านเซิร์ฟเวอร์ตัวเดียวกัน โครงสร้างแบบสตาร์จึงง่ายต่อการติดตั้งและการบำรุงรักษา ข้อเสียของรูปแบบสตาร์คือ ถ้าซีแอนซีเซิร์ฟเวอร์ไม่สามารถใช้งานได้ด้วยเหตุผลใดๆ ก็ตาม บอทจะยังสามารถทำงานได้บนเครื่องที่ถูกคุกคาม แต่จะไม่สามารถรับคำสั่งใดๆ และไม่สามารถส่งข้อมูลไปยังบอทมาสเตอร์ได้



### 2) มัลติเซิร์ฟเวอร์ (Multi-Server)

โครงสร้างมัลติเซิร์ฟเวอร์คล้ายกับโครงสร้างแบบสตาร์ แต่โครงสร้างมัลติเซิร์ฟเวอร์จะมีเซิร์ฟเวอร์อย่างน้อย 2 ตัวที่บอทนั้นสื่อสารด้วยแสดงดังรูปที่ 2.3 โครงสร้างแบบนี้จะลดปัญหาที่เกิดกับโครงสร้างแบบสตาร์ได้ เนื่องจากถ้าเซิร์ฟเวอร์หนึ่งไม่สามารถทำงานได้ด้วยเหตุผลใดก็ตาม เซิร์ฟเวอร์อื่นก็จะสามารถควบคุมบอทเหล่านั้นแทนได้

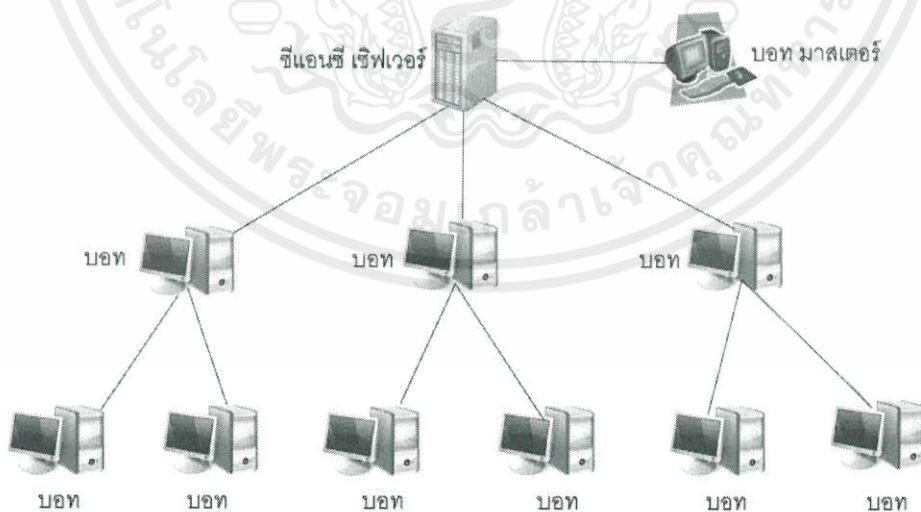
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 2.3 ตัวอย่างโครงสร้างแบบมีลติเซิร์ฟเวอร์ของซีแอนซีบอทเน็ต

### 3) ลำดับชั้น (Hierarchical)

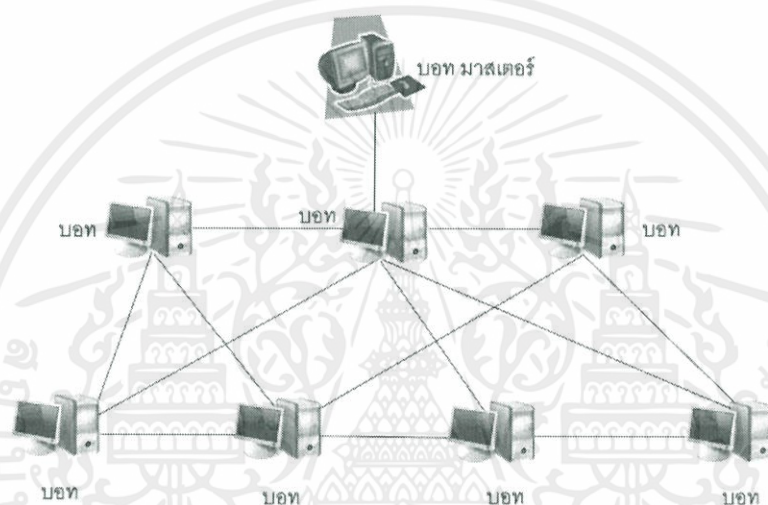
โครงสร้างแบบลำดับชั้นจะใช้โครงสร้างย่อยของผู้ที่อยู่ใต้บังคับบัญชาแทนการส่งการผ่านซีแอนซีเซิร์ฟเวอร์เพียงตัวเดียว มีรูปแบบคล้ายทรี สามารถกระจายคำสั่งจากบอทตัวหนึ่งไปยังบอทตัวอื่นๆ ได้ดังรูปที่ 2.4 ข้อเสียของโครงสร้างแบบลำดับชั้น คือ บอทเพียงตัวเดียวไม่เพียงพอที่จะทราบสถานที่และจำนวนของบอททั้งหมดที่มีได้



เอกสารนี้เป็นเอกสารที่รูปที่ 2.4 ตัวอย่างโครงสร้างแบบลำดับชั้นของซีแอนซีบอทเน็ตไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

#### 4) สุ่ม (Random)

โครงสร้างแบบสุ่มเป็นโครงสร้างที่ไม่มีศูนย์กลางแต่จะติดต่อสื่อสารโดยตรงผ่านเซิร์ฟเวอร์ ดังแสดงในรูปที่ 2.5 ตัวอย่างนี้ของโครงสร้างแบบสุ่ม เช่น พีทูพี peer-to-peer (P2P) ซึ่งคำสั่งใน พีทูพี นั้นจะมีการแนะนำบอทในบอทเน็ตและสามารถแพร่กระจายไปยังบอททุกตัวแทนบอทตัวอื่นได้ โครงสร้างแบบสุ่มนี้ยากที่จะปิดตัวลง เนื่องจากไม่มีศูนย์กลางของเซิร์ฟเวอร์ซีแอนซี และมีข้อเสียคือการสื่อสารอาจล่าช้าและไม่สามารถคาดเดาได้



รูปที่ 2.5 ตัวอย่างโครงสร้างแบบสุ่มของซีแอนซีบอทเน็ต

#### 2.2.4 โพรโทคอลในการสื่อสาร

การติดต่อสื่อสารกันระหว่างบอทเน็ตจะขึ้นอยู่กับโพรโทคอลที่ถูกนำมาใช้ รูปแบบและโครงสร้างของการสื่อสาร จะถูกใช้เพื่อส่งข้อมูลไปยังซีแอนซีเซิร์ฟเวอร์ และรอรับคำสั่งจากบอทมาสเตอร์ แต่เดิมมีการสื่อสารผ่านทางโพรโทคอลไออาร์ซี (Internet Relay Chat) โดยบอทจะเข้าไปทำงานร่วมกับไออาร์ซีเซิร์ฟเวอร์ที่มีการป้องกันแบบทั่วไปด้วยรหัสผ่าน บอทจะดักรับคำสั่งผ่านทางแชนแนลไออาร์ซี แต่การใช้งานเอสเอสแอล (SSL) ในไออาร์ซี เพื่อเข้ารหัสข้อมูลและคำสั่ง อาจสามารถหลีกเลี่ยงการตรวจจับของบอทได้ แต่ในปัจจุบันบอทส่วนมากใช้งานโพรโทคอลเอชทีทีพี (HTTP) เพราะไฟร์วอลล์นั้นมีการบล็อกทราฟฟิกที่มีการติดต่อกับซีแอนซีเซิร์ฟเวอร์ แต่เครือข่ายเกือบทั้งหมดจะอนุญาตให้โพรโทคอลเอชทีทีพีผ่านไปได้ เพราะใช้ในการเข้าถึงเว็บไซต์ และยิ่งไปกว่านั้น บางส่วนของบอทเน็ตมีการนำโพรโทคอลพีทูพี (Peer to Peer) มาใช้ในโครงสร้างแบบสุ่ม ทำให้ยากต่อการตรวจจับ และยากที่จะปิดตัวลงอีกด้วย ตัวอย่างบอทจำพวกพีทูพี ได้แก่ สพายบอท

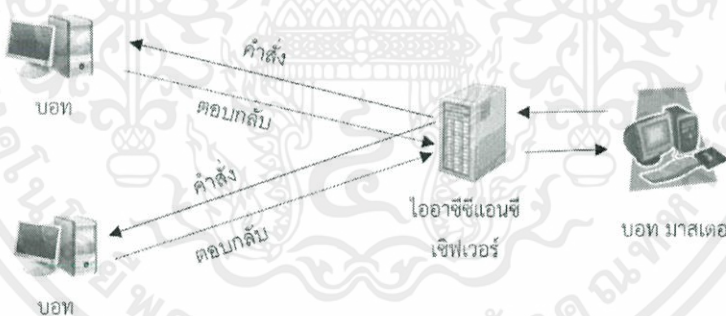
(SpyBot) และ อโกบอท (Agobot) อย่างไรก็ตามการใช้งานโพรโทคอลไออาร์ซีก็ยังเป็นที่นิยมมากที่สุด เนื่องจากมีสคริปต์ที่พร้อมใช้งาน และมีทรัพยากรวัสดุต่างๆ มากมาย

## 2.2.5 กลไกการควบคุมและสั่งการศูนย์กลาง

ซีแอนซีเซิร์ฟเวอร์ สามารถแบ่งประเภทโดยอยู่บนพื้นฐานของโพรโทคอลที่ใช้ในการสื่อสารและความสามารถในการสร้างการเชื่อมต่อระหว่างกันได้เป็น 2 ประเภท คือ ไออาร์ซี และ เอชทีทีพี

### 1) ไออาร์ซีบอทเน็ต (IRC-base Botnets)

ไออาร์ซีเป็นระบบที่ถูกใช้โดยผู้ใช้คอมพิวเตอร์ในการสื่อสารออนไลน์หรือพูดคุยกันแบบเรียลไทม์ วิธีนี้ถูกนำมาใช้ในบอทรุ่นแรก โดยบอทมาสเตอร์จะใช้ไออาร์ซีเซิร์ฟเวอร์และเซิร์ฟเวอร์ที่เกี่ยวข้องเพื่อกระจายคำสั่ง แต่ละบอทจะต้องติดต่อไปที่ไออาร์ซีเซิร์ฟเวอร์และเซิร์ฟเวอร์ที่บอทมาสเตอร์ได้เลือกและจะรอคำสั่งจากบอทมาสเตอร์ที่เซิร์ฟเวอร์นี้ บอทมาสเตอร์จะสร้างการสื่อสารแบบเรียลไทม์ในทุกการเชื่อมต่อกับบอท และไออาร์ซีบอทจะคอยส่งคำสั่งอย่างเดียวซึ่งหมายความว่าเมื่อไออาร์ซีบอทเชื่อมต่อไปที่เซิร์ฟเวอร์ที่ได้เลือกไว้ บอทจะไม่ถูกตัดการเชื่อมต่อและยังคงสภาพในโหมดเชื่อมต่อตลอดเวลา ดังรูปที่ 2.6

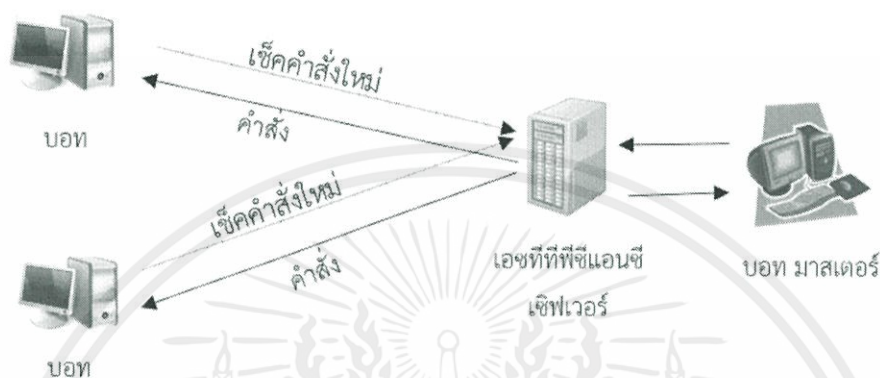


รูปที่ 2.6 ไออาร์ซีเบสซีแอนซีบอทเน็ต

### 2) เอชทีทีพีบอทเน็ต (HTTP-base Botnets)

เอชทีทีพีซีแอนซีเป็นเทคโนโลยีใหม่ที่ช่วยให้บอทมาสเตอร์สามารถควบคุมบอทโดยใช้เอชทีทีพีโพรโทคอล ในเทคนิคนี้บอทจะใช้ยูอาร์แอลที่เจาะจงหรือไอพีแอดเดรสที่กำหนดโดยบอทมาสเตอร์ เพื่อเชื่อมต่อไปยังเว็บเซิร์ฟเวอร์เฉพาะ เอชทีทีพีบอทจะใช้วิธีการดึง (Pull) ซึ่งต่างจากวิธีการผลัก (Push) ของไออาร์ซีบอท ในวิธีการดึงนั้นเอชทีทีพีบอทจะไม่คงสถานะอยู่ในโหมดการเชื่อมต่อหลังจากที่สร้างการเชื่อมต่อไปที่ซีแอนซีเซิร์ฟเวอร์แล้ว ในขั้น

แรกในวิธีการดึง บอทมาสเตอร์จะเผยแพร่คำสั่งที่บางเซิร์ฟเวอร์ และบอทจะเข้าไปที่เว็บไซต์ของบอทเพื่ออัปเดตตัวเองหรือรับคำสั่งใหม่เป็นระยะๆ ด้วยกระบวนการที่ต่อเนื่องในช่วงเวลาปกติ ซึ่งจะถูกกำหนดโดยบอทมาสเตอร์



รูปที่ 2.7 เอชทีทีพีเบสซีแอนซีบอทเน็ต

### 2.3 เหตุผลที่เลือกเอชทีทีพีบอทเน็ตในการวิจัย

ในรุ่นแรกของบอทเน็ตเทคโนโลยีไออาร์ซีถูกใช้โดยบอทมาสเตอร์เพื่อควบคุมบอทเพราะระบบไออาร์ซีมีข้อได้เปรียบหลายอย่าง เช่น การใช้งานที่สะดวกทั้งการควบคุม,การจัดการ อย่างไรก็ตามจุดอ่อนของไออาร์ซีบอทเน็ตคือกลไกการควบคุมศูนย์กลางบอทเน็ตทั้งหมดสามารถถูกทำลายได้โดยการบล็อกไออาร์ซีเซิร์ฟเวอร์หรือบล็อกพอร์ตไออาร์ซี ดังนั้นบอทแบบพีทูพีจึงถูกออกแบบมาเพื่อแก้ปัญหานี้ในบอทแบบพีทูพีจะไม่มีศูนย์กลางตัวควบคุมและสั่งการ แต่จะมีการกระจายเซิร์ฟเวอร์คำสั่งจะถูกส่งไปหาบอทโดยบอทเอง นอกจากนี้บางวิธีการถอดรหัสถูกนำมาใช้เพื่อการสื่อสารที่ปลอดภัยเทคนิคนี้ทำให้ยากขึ้นสำหรับตรวจจับพีทูพีบอทเน็ตเมื่อเทียบกับไออาร์ซีบอทเน็ต อย่างไรก็ตามพีทูพีบอทเน็ตไม่ได้ถูกนำมาใช้อย่างกว้างขวางอย่างไออาร์ซีบอทเน็ตเพราะการดำเนินการและการควบคุมบอทเน็ตของพีทูพีนั้นค่อนข้างยากและซับซ้อน นอกจากนี้ยังมีความล่าช้าในการส่งคำสั่งระหว่างพีทูพีบอทด้วยกันและบอทมาสเตอร์ยังไม่สามารถที่จะรู้เกี่ยวกับสถานะของคำสั่งด้วย

ในปัจจุบันบอทมาสเตอร์ได้เริ่มใช้ศูนย์กลางเป็นซีแอนซีอีกครั้ง เอชทีทีพีโพรโทคอลถูกใช้แทนที่ไออาร์ซีโพรโทคอลและใช้พอร์ต 80 เพราะว่าความหลากหลายของการบริการ ทำให้ยากที่จะบล็อกศูนย์กลางซีแอนซีเซิร์ฟเวอร์นี้ นอกจากนี้ในการใช้เอชทีทีพีโพรโทคอลบอทจะซ่อนการสื่อสารผ่านทางกระแสของเอชทีทีพีปกติและจะหลีกเลี่ยงการตรวจจับโดยตัวตรวจจับเน็ตเวิร์ค อย่างเช่นไฟร์วอลล์ เป็นต้น

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ในการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากการทบทวนในลักษณะของ ไออาร์ซี,พีทูพี และ เอชทีทีพีบอทเน็ตแล้วเป็นที่ชัดเจนว่า ซี แอนซีของเอชทีทีพีนั้นเป็นเทคโนโลยีที่ใหม่ที่เป็นที่ต้องการของบอทเน็ตเมื่อเทียบกับไออาร์ซีและพีทูพีบอทเน็ตแล้ว เอชทีทีพีเบสบอทเน็ตจะมีเซตของคุณสมบัติที่ทำให้เป็นเรื่องยากในการตรวจจับแต่เป็นเรื่องน่าแปลกใจที่จำนวนของงานวิจัยที่มุ่งเน้นไปที่การตรวจสอบของบอทเน็ตที่ใช้ เอชทีทีพี นั้นค่อนข้างต่ำเมื่อเทียบกับจำนวนของงานวิจัยเกี่ยวกับวิธีการตรวจสอบสำหรับบอทเน็ตที่ที่ใช้ ไออาร์ซีหรือพีทูพี

## 2.4 วิธีการที่ใช้ในการตรวจสอบบอทเน็ต

วิธีการที่ใช้ในการตรวจสอบบอทเน็ตมีดังต่อไปนี้

### 2.4.1 ฮันนีพอทและฮันนีเน็ต

ฮันนีพอทและฮันนีเน็ต (Honeypot และ Honeynet) เป็นเครื่องมือที่ใช้เป็นกับดักสำหรับดักจับบอทที่สามารถตรวจจับได้ หรือรวบรวมข้อมูลของกิจกรรมของบอทได้ ข้อมูลที่สามารถนำมาใช้เพื่อทำความเข้าใจเกี่ยวกับพฤติกรรมของบอทหรือเจตนาของบอทมาสเตอร์ จะใช้วิธีการรวบรวมไบนารีโค้ดของบอทและข้อมูลอื่นๆ ที่เกี่ยวกับบอท

### 2.4.2 ตรวจจับโดยลายเซ็น

ลายเซ็นหรือซิกเนเจอร์ (Signature) หมายถึงรูปแบบที่รู้จักกันหรือลักษณะของภัยคุกคามจากผู้บุกรุกเข้าสู่ระบบคอมพิวเตอร์ โดยจะวิเคราะห์และเปรียบเทียบรูปแบบเหล่านี้หรือลักษณะเหล่านี้ ซึ่งมีความเป็นไปได้ที่จะแยกแยะกิจกรรมที่เป็นภัยคุกคามจากกิจกรรมทั่วไป

การตรวจจับแบบวิธีลายเซ็นนั้นเป็นวิธีที่ประสิทธิภาพไม่เท่าที่ควร เนื่องจากวิธีนี้ไม่สามารถระบุรูปแบบพฤติกรรมหรือลักษณะใหม่ๆ ของบอทได้ วิธีการนี้จะขึ้นอยู่กับเปรียบเทียบข้อมูลที่เก็บรวบรวมไว้สำหรับที่รูปแบบที่รู้จักอยู่แล้ว ดังนั้นวิธีการนี้เป็นวิธีการที่ดีที่สุดสำหรับตรวจสอบบอทที่รู้จักกันดี แต่ค่อนข้างไร้ประโยชน์สำหรับการตรวจสอบบอทใหม่

### 2.4.3 ตรวจจับโดยการตรวจสอบดีเอ็นเอส

การตรวจสอบและวิเคราะห์กราฟฟิคด้วยดีเอ็นเอส (DNS Monitoring) ถูกนำมาใช้เป็นเทคนิคในการตรวจจับบอท โดยพบว่าบอทจะสร้างกราฟฟิคดีเอ็นเอสในบางสถานการณ์ เช่น เมื่อมีการระบุเซิร์ฟเวอร์ซีแอนซีหรือมีการเตรียมการโจมตี เช่น ดีดีไอเอส เป็นต้น

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## 2.5 การตรวจจับโดยการวิเคราะห์พฤติกรรมของเครือข่าย

การตรวจจับโดยการวิเคราะห์พฤติกรรมของเครือข่าย (Detection Based on Network Behavior Analysis) เรียกว่า เอนบีเอ (NBA) เป็นวิธีที่สามารถใช้ในการเก็บรวบรวมความหลากหลายของข้อมูล และสถิติเกี่ยวกับทราฟฟิกของเครือข่าย จากนั้นจะนำข้อมูลมาวิเคราะห์เพื่อตรวจจับสัญญาณที่เป็นภัยคุกคามหรือกิจกรรมที่เป็นอันตรายต่างๆ วิธีเอนบีเอประกอบไปด้วยองค์ประกอบที่หลากหลายรวมถึงเซ็นเซอร์และตัวจัดการเซิร์ฟเวอร์ด้วย โดยระบบเอนบีเอนั้นจะรวบรวมข้อมูล เช่น ไอพีแอดเดรส, ระบบปฏิบัติการ, การบริการที่เปิดให้ใช้บริการ และข้อมูลที่ถูกบันทึกไว้ เช่น Timestamp, ประเภทของเหตุการณ์ (event type), โพรโทคอลเครือข่าย, พอร์ตโฮส และ ฟิลด์ของเฮดเดอร์แพ็คเก็ตสำหรับแต่ละโคลเอนต์

### 2.5.1 ตัวอย่างการตรวจจับโดยการวิเคราะห์พฤติกรรมตามช่วงของเวลา (Period)

พฤติกรรมของเอชทีทีพีเบสโตนีต (HTTP-base Botnets) จะมีการเริ่มการเชื่อมต่อและร้องขอเอชทีทีพีเก็ท (HTTP GET) ไปยังเว็บไซต์ เพื่อร้องขอคำสั่ง หลังจากที่เว็บไซต์ของซีแอนซีได้รับคำร้องขอแล้วก็จะส่งคำสั่งกลับมาพร้อมกับเครื่องคอมพิวเตอร์ที่ติดบอท ข้อดีของวิธีนี้ คือ เว็บไซต์ซีแอนซีไม่จำเป็นต้องติดแท็ก (Track) และเครื่องคอมพิวเตอร์ที่ติดบอทเพียงแค่อุปกรณ์คอมพิวเตอร์ที่ติดบอทด้วยกันติดต่อกันมาแทน ข้อเสียของวิธีนี้ คือ เว็บไซต์นั้นไม่สามารถควบคุมหรือรักษาการเชื่อมต่อของเครื่องคอมพิวเตอร์ที่ติดบอทได้จนกว่าเครื่องนั้นๆ จะติดต่อกลับมาที่เซิร์ฟเวอร์เอง

กลไกที่ทำให้เครื่องคอมพิวเตอร์ติดบอท คือ การฝังโค้ดที่เป็นอันตรายลงไป และมีการเรียกใช้คำสั่ง sleep() (คำสั่งที่ใช้ในการรอตามระยะเวลาที่กำหนดก่อนจะมารับคำสั่งใหม่ที่เว็บไซต์ซีแอนซี) ใน while loop ซึ่งแสดงได้ดังรูปที่ 2.8 แสดงตัวอย่างโค้ดที่อาจจะถูกนำมาใช้ และตารางที่ 2.1 แสดงตารางบันทึกเวลา (Timestamps) การเชื่อมต่อแอบเอชทีทีพีของแต่ละการเชื่อมต่อ

```
while (1) {
    get_and_process_cnc_commands();
    sleep(60);
}
...
void get_and_process_cnc_commands()
{
    // Insert code here to connect to C&C website, poll for a command
    // and process the command
}
```

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งรูปที่ 2.8 โค้ดแสดงกลไกการใช้คำสั่ง sleep() อย่างง่าย

ตารางที่ 2.1 ตารางเวลาในแต่ละช่วงเวลาที่ยอดที่มีการติดต่อกับซีแอนซีเซิร์ฟเวอร์

Date/Time	Time Interval (Seconds)
2010-05-01 21:47:46	-
2010-05-01 22:47:50	3604
2010-05-01 23:47:53	3603
2010-05-01 00:47:55	3602
2010-05-01 01:47:58	3603
2010-05-01 02:47:01	3603
2010-05-01 03:47:04	3603
2010-05-01 04:47:07	3603
2010-05-01 05:47:10	3603
2010-05-01 06:47:13	3603
2010-05-01 07:47:16	3603
2010-05-01 08:47:19	3603
2010-05-01 09:47:22	3603
mean = 3603 std. dev = 0.43	

จากการสังเกตเครื่องคอมพิวเตอร์ที่ติดบอหจะมีพฤติกรรมในการติดต่อกลับมาที่ซีแอนซีเซิร์ฟเวอร์เป็นระยะๆ ทำให้สามารถนำพฤติกรรมการติดต่อดังกล่าว ซึ่งมีช่วงระยะเวลาเกือบจะคงที่ มาใช้ในการคำนวณและตรวจจับการเชื่อมต่อที่มีระยะเวลา

### 2.5.2 เหตุผลที่เลือกวิธีการตรวจจับโดยใช้การวิเคราะห์พฤติกรรมเครือข่าย

- ความสามารถในการตรวจจับภัยคุกคามที่ไม่รู้จัก

บอทมาสเตอร์จะทำการอัปเดตข้อมูลวันต่อวัน เพื่อหลีกเลี่ยงการตรวจจับ แต่ระบบเอนบีเอนั้นสามารถขัดขวางกลยุทธ์ต่างๆ ของบอทมาสเตอร์เหล่านี้ได้ และด้วยคุณสมบัติของเอนบีเอนี้ทำให้สามารถเพิ่มประสิทธิภาพในการตรวจจับบอทเน็ตได้

- ความสามารถในการตรวจจับภัยคุกคามที่มีการเข้ารหัส

บอทเน็ตจะพยายามซ่อนทราฟฟิกของบอทกับทราฟฟิกทั่วไป หรือใช้วิธีการเข้ารหัส เอกสารนี้เป็นเอกสารที่เอนบีเอนีสามารถมองรูปแบบของการไหลของข้อมูลในเครือข่ายที่ผิดปกติได้ นอกจากนี้เอนบีเอนียังสามารถวิเคราะห์ข้อมูลจราจรที่เข้ารหัสได้ ซึ่งช่วยให้การตรวจจับแบบเอนบีเอนีสามารถให้ผลลัพธ์ที่ดีที่สุดเมื่อมีการทำงานร่วมกับวิธีอื่น

## 2.6 เทคนิคการแบ่งกลุ่มแบบเชิงลำดับชั้น

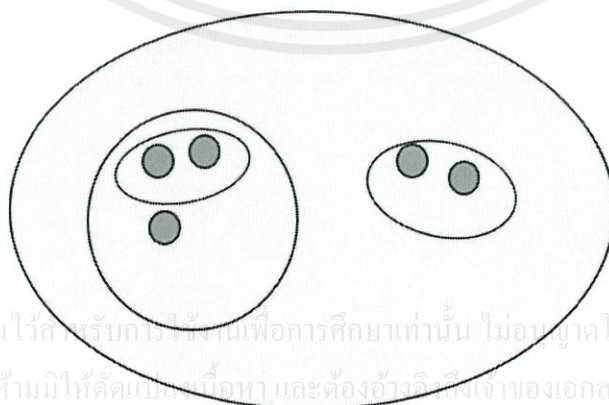
เทคนิคการแบ่งกลุ่มแบบเชิงลำดับชั้น (Hierarchical Cluster) มีเครื่องมือในการวิเคราะห์อย่างแพร่หลายเป็นจำนวนมาก แนวคิด คือ การสร้างบาลานซ์ทรีของข้อมูลที่เป็นลำดับ และรวมกลุ่มของจุดที่คล้ายกัน เพื่อสรุปการใช้งานของข้อมูล ข้อดีของเทคนิคแบบเชิงลำดับชั้นเมื่อเปรียบเทียบกับเทคนิคแบบเคมีน คือ ข้อมูลจำเป็นสำหรับเทคนิคแบบเคมีน (K-mean) คือ จำนวนของกลุ่ม, การมอบหมายข้อมูลกับกลุ่ม, ระยะทางระหว่างข้อมูล แต่ข้อมูลสำหรับเทคนิคแบบเชิงลำดับชั้นต้องการเพียงตัวชี้วัดของความคล้ายคลึงกันระหว่างจุดของกลุ่มข้อมูล ในงานวิจัยนี้จึงได้เลือก เบิร์ชอัลกอริทึม (BIRCH Algorithm) ในการแบ่งกลุ่มของข้อมูล

### 2.6.1 ประโยชน์ของอัลกอริทึมเบิร์ช

- 1) อัลกอริทึมเบิร์ชอาศัยการจัดกลุ่มโดยการใช้ข้อมูลที่ดินมีอยู่ โดยไม่จำเป็นต้องสแกนข้อมูลทั้งหมดของระบบใหม่อีกครั้ง
- 2) อัลกอริทึมเบิร์ชสามารถสร้างกลุ่มของข้อมูลได้โดยอาศัยการสแกนระบบทั้งหมดเพียงครั้งเดียว
- 3) อัลกอริทึมเบิร์ชสามารถบริหารจัดการการใช้หน่วยความจำสำหรับการแบ่งกลุ่มข้อมูลได้

### 2.6.2 หลักการของอัลกอริทึมเบิร์ช

อัลกอริทึมเบิร์ชจะเริ่มแบ่งกลุ่มของข้อมูลที่กลุ่มของจุดๆ หนึ่ง (ทุกจุดในฐานะข้อมูลในกลุ่ม) หลังจากนั้นภายในแต่ละกลุ่มจะหาจุดที่ใกล้เคียงที่สุด แล้วรวมกันเป็นอีกกลุ่มหนึ่ง และทำแบบนี้ไปเรื่อยๆ จนกระทั่งเหลือเพียงกลุ่มใหญ่สุดกลุ่มเดียว การคำนวณของกลุ่มจะทำได้ด้วยระยะทางเมตริก ( $O(n^2)$ ) และ เวลา  $O(n^2)$



รูปที่ 2.9 การแบ่งกลุ่มแบบเชิงลำดับชั้น

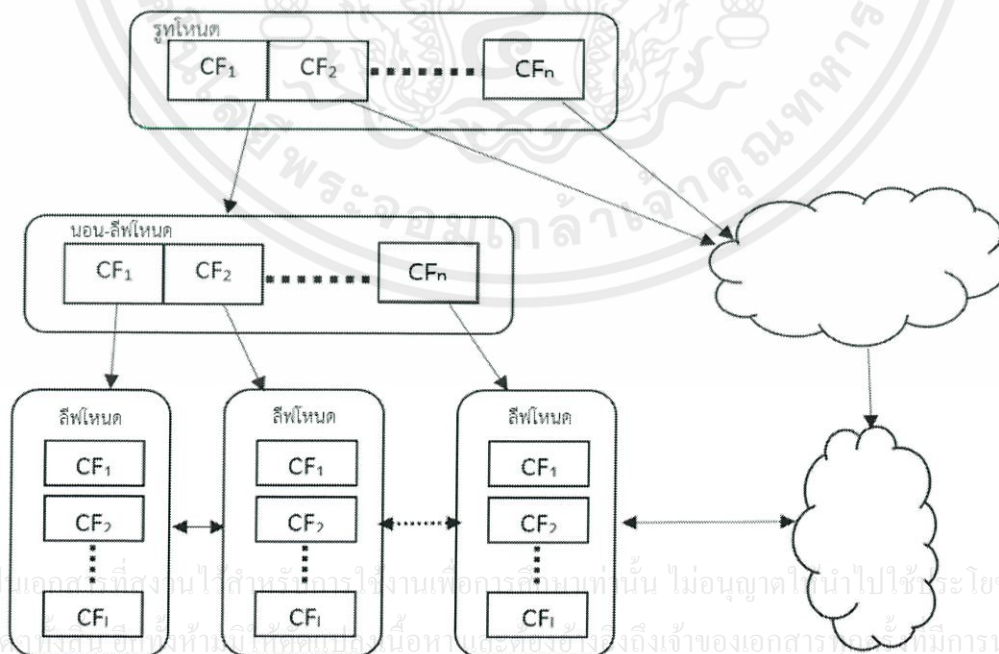
### 2.6.2.1 ลักษณะการแบ่งกลุ่ม

อัลกอริทึมเบียร์จะสร้างซีเอฟทรี (Clustering Feature Tree) ในขณะที่มีการสแกนชุดข้อมูล ซึ่งแต่ละเอนทรีในซีเอฟทรีจะแสดงกลุ่มของออปเจ็ค และลักษณะที่โดดเด่น 3 ประการ (N,LS,SS) โดยกำหนดให้ N เป็นเวกเตอร์ของจุดข้อมูลในกลุ่มที่มี d-มิติ และกำหนดให้  $X_i$  ( $i = 1, 2, 3, \dots, N$ ) เป็นซีเอฟเวกเตอร์ของกลุ่ม รูปแบบของ CF เป็น (N,LS,SS) โดยที่ N คือ จำนวนจุดของข้อมูลในกลุ่ม, LS คือ ผลรวมเชิงเส้นของ N จุดข้อมูล และ SS คือ ผลรวมตารางของ N จุดข้อมูล

### 2.6.2.2 ซีเอฟทรี

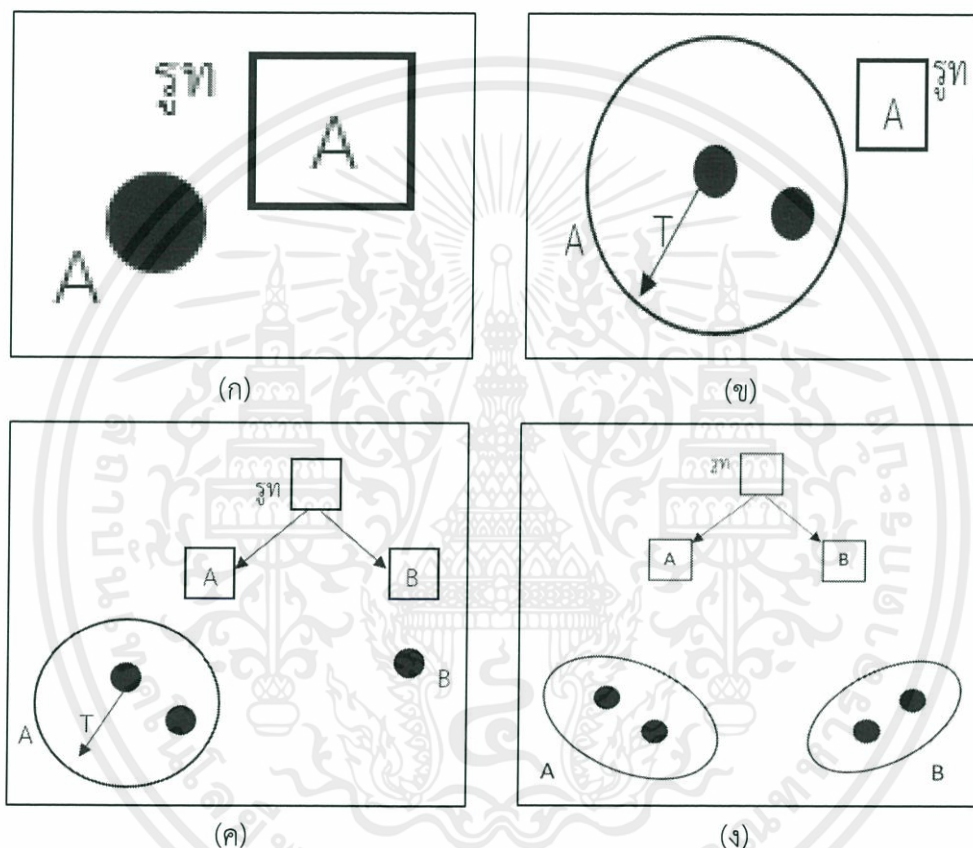
ความสูงของบาลานซ์บริจจะมี 2 ค่า คือ บรานซ์เฟคเตอร์ B และ เทรชโฮล T ซึ่งแต่ละนอนลิฟโหนดจะมีเอนทรีได้มากที่สุด B เอนทรี ในรูป  $[CF_i, child_i]$  โดยที่  $child_i$  เป็นจุดที่ชี้ไปยังโหนดที่ i และ  $CF_i$  เป็นซีเอฟของกลุ่มย่อยที่ถูกแทนด้วย  $child_i$  นี้ ดังนั้น นอนลิฟโหนดจะแสดงให้เห็นถึงกลุ่มที่สร้างขึ้นจากกลุ่มย่อยทั้งหมดที่แสดงถึงเอนทรีของตัวเอง

ลิฟโหนดจะมีเอนทรีมากที่สุด L เอนทรี แต่ละโหนดจะมีรูปแบบ  $[CF_i]$  โดยที่  $i = 1, 2, \dots, L$  นอกจากนี้ยังมีตัวซ้อยู่ 2 ตัวคือ ก่อนหน้าและถัดไป ซึ่งถูกใช้ในการเชื่อมลิฟโหนดเข้าด้วยกันเพื่อให้เวลาสแกนมีประสิทธิภาพ ลิฟโหนดยังแสดงให้เห็นถึงกลุ่มที่สร้างมาจากกลุ่มย่อยทั้งหมดที่แสดงถึงเอนทรีตัวเอง แต่ทุกเอนทรีในลิฟโหนดต้องมีการตอบสนองค่าเทรชโฮลด้วยการพิจารณาถึงค่าเทรชโฮล T ด้วย โดยที่ค่ากลางนั้นต้องมีค่าน้อยกว่า T



รูปที่ 2.10 ซีเอฟทรีของแต่ละโหนดย่อย

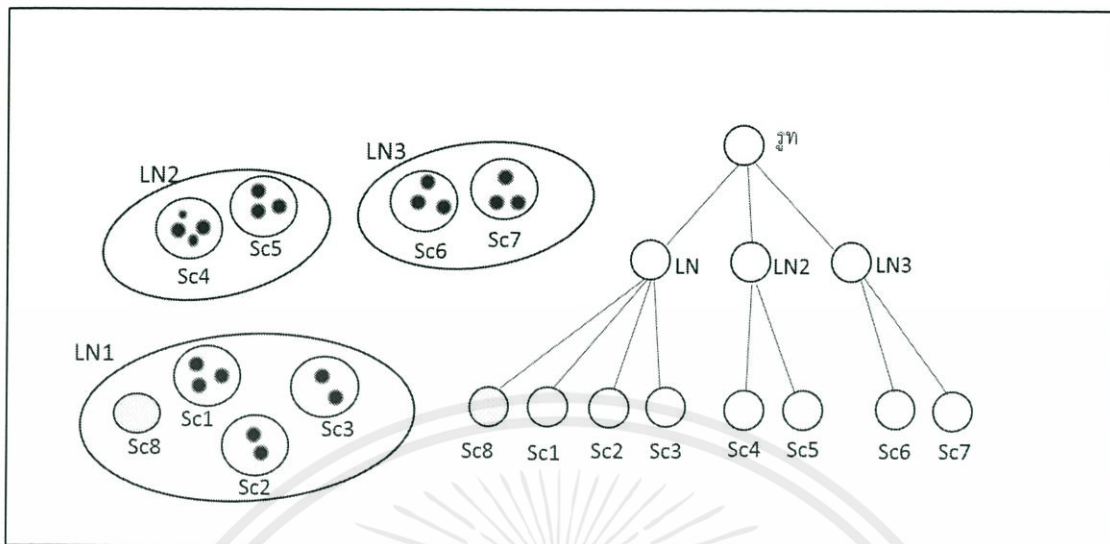
ตัวอย่างของซีเอฟที ในช่วงเริ่มต้นของจุดข้อมูลในกลุ่มหนึ่ง ดังรูปที่ 2.11(ก) เมื่อมีข้อมูลใหม่เข้ามาจะทำการตรวจสอบว่าระยะห่างนั้นมีค่ามากกว่า T หรือไม่ ดังรูปที่ 2.11(ข) ถ้าขนาดของกลุ่มนั้นมีขนาดใหญ่เกินกว่าค่า T กลุ่มนั้นจะถูกแบ่งออกเป็น 2 กลุ่ม และจะกระจายจุดออกจากกัน ดังรูปที่ 2.11(ค) ที่แต่ละโหนดของทรี ซีเอฟทีจะเก็บข้อมูลเกี่ยวกับค่าเฉลี่ยของกลุ่มและเฉลี่ยผลรวมที่ได้ เพื่อให้สามารถคำนวณขนาดของกลุ่มได้อย่างมีประสิทธิภาพ ดังรูปที่ 2.11(ง)



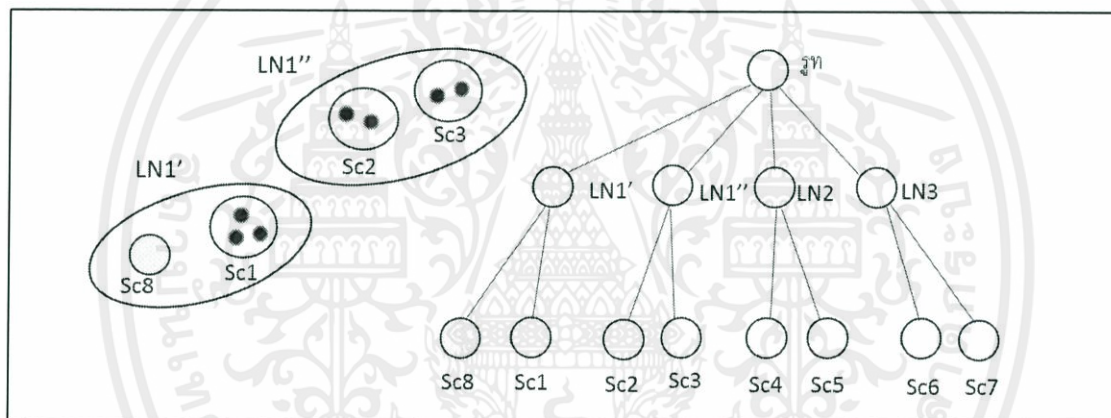
รูปที่ 2.11 ตัวอย่างซีเอฟที

ตัวอย่างอื่นๆ ของซีเอฟที แสดงในรูปที่ 2.12 จะเห็นได้ว่ามีข้อมูล sc8 เข้ามาทำให้กลุ่ม LN1 มีมากกว่าสมาชิก 3 จึงต้องทำการแบ่งกลุ่ม LN1 ออกเป็น 2 กลุ่ม ดังรูปที่ 2.12(ค)

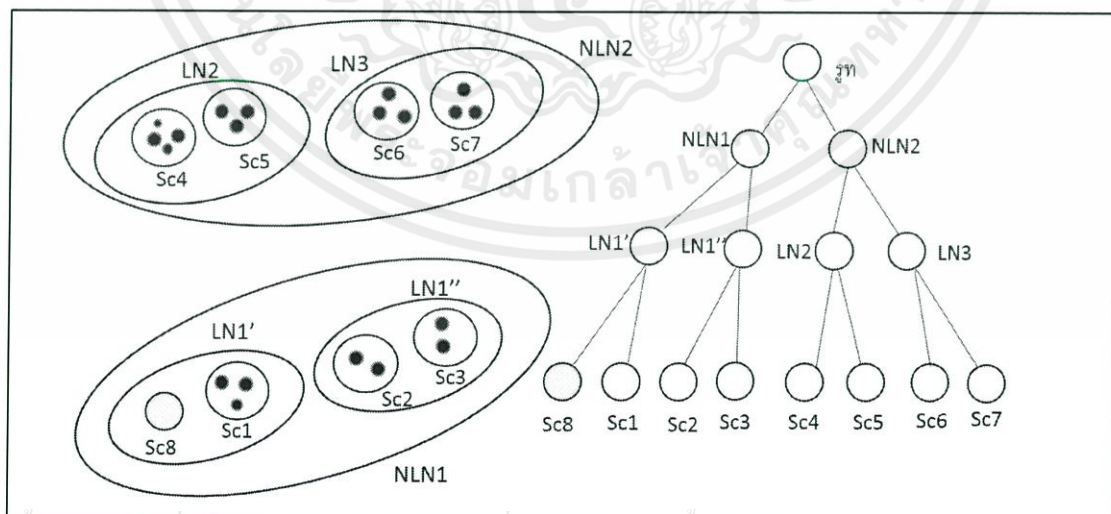
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



(ก)



(ข)



(ค)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดลอกเผยแพร่ข้อมูลและข้อมูลเชิงลึกของเอกสารทุกครั้งที่มีการนำไปใช้

รูปที่ 2.12 ตัวอย่างอื่นๆของซีเอฟทีรี

จากรูปที่ 2.12(ข) จะเห็นว่าอนาลิฟโหนดนั้นมีข้อมูลมากกว่า 3 จึงต้องให้ Root นั้นแบ่งกลุ่มย่อยออกไปเป็นดังรูปที่ 2.11(ค)

สรุปขั้นตอนของอัลกอริทึมเบิร์ชได้ดังนี้

- ขั้นตอนที่ 1 : สแกนข้อมูลทั้งหมดและสร้างจุดเริ่มต้นในหน่วยความจำในซีเอฟพี
- ขั้นตอนที่ 2 : ทำการสแกนลีฟโหนดเพื่อสร้างทรีใหม่ที่ให้มีขนาดเล็กลง
- ขั้นตอนที่ 3 : จัดกลุ่มที่ใหญ่ที่สุด
  - จะใช้ Global หรือ Semi-Global อัลกอริทึมเพื่อแบ่งกลุ่มโหนดลีฟโหนดทั้งหมด
  - Adapted agglomerative hierarchical clustering algorithm ถูกนำไปใช้โดยตรงกับกลุ่มย่อยที่แสดงโดยเวกเตอร์ซีเอฟพีนั้นๆ
- ขั้นตอนที่ 4 : ปรับปรุงผลลัพธ์การแบ่งกลุ่ม (ทางเลือก)
  - เพื่อตรวจสอบข้อมูลและแก้ไขข้อมูลที่ไม่ถูกต้องและปรับกลุ่มต่อไป
  - ใช้จุดศูนย์กลางของกลุ่มในการสร้างที่ขั้นตอนที่ 3 เป็น Seed และ รีดิสทริบิวต์ข้อมูลไปยัง Seed ใกล้เคียงเพื่อเซตกลุ่มข้อมูลใหม่

สรุปอัลกอริทึมเบิร์ช

- อัลกอริทึมเบิร์ชจะดำเนินการได้รวดเร็วกว่าอัลกอริทึมอื่นๆ ที่มีอยู่ และสามารถใช้งานกับชุดข้อมูลที่มีขนาดใหญ่ได้
- อัลกอริทึมเบิร์ชจะสแกนข้อมูลทั้งหมดภายในครั้งเดียว
- สามารถตรวจจับค่าผิดปกติได้ดี
- เมื่อเปรียบเทียบกับอัลกอริทึมอื่นพบว่าอัลกอริทึมเบิร์ชมีความเสถียรและยืดหยุ่นกว่า

## 2.7 ระยะห่างเชิงเลขคณิต, ระยะห่างเชิงประเภท และระยะห่างเชิงลำดับชั้นสำหรับอัลกอริทึมเบิร์ช

หัวข้อนี้จะกล่าวถึงระยะห่างเชิงเลขคณิต, ระยะห่างเชิงประเภท และระยะห่างเชิงลำดับชั้นสำหรับอัลกอริทึมเบิร์ช (Numerical, Categorical and Hierarchical Distance with BIRCH Algorithm) เนื่องจากปัญหาสำคัญอย่างหนึ่งในการนำข้อมูลของการใช้งานเครือข่ายมาคำนวณทางคณิตศาสตร์คือข้อมูลที่ได้มาจากการเก็บรวบรวมจากอุปกรณ์เครือข่ายนั้นส่วนมากมักจะเป็นข้อมูลที่ไม่สามารถนำมาคำนวณได้โดยตรงยกตัวอย่างข้อมูล เช่น หมายเลขไอพีต้นทาง/ปลายทาง, หมายเลขพอร์ตต้นทาง/ปลายทาง หรือ ชนิดของโพรโทคอล เป็นต้น การนำข้อมูลหรือคุณสมบัติของแพ็คเก็ต

ดังกล่าวมาวิเคราะห์หรือคำนวณทางคณิตศาสตร์จำเป็นต้องมีการแทนค่าข้อมูลนั้นๆ และกำหนดวิธีการคำนวณอย่างเหมาะสมแทนการคำนวณแบบปกติ จากการศึกษาพบว่าข้อมูลดังกล่าวนี้สามารถถูกแทนค่าและกำหนดการดำเนินการทางคณิตศาสตร์ที่เหมาะสมให้ได้

กำหนดให้  $R$  แทนเซตของจำนวนจริง (Real Number),  $Z$  แทนเซตของจำนวนเต็ม (Integers),  $R^d$  แทนเวกเตอร์ซึ่งมี  $d$  มิติบนเซตของ  $R$ ,  $C$  แทนกลุ่มของคุณสมบัติ (Cluster) ซึ่งมีจำนวนสมาชิกในกลุ่มเป็น  $N$  เขียนได้เป็น  $C = \{X_j : j = 1, \dots, N\}$ , โดยที่  $X$  เป็นเวกเตอร์ซึ่งมี  $d$  มิติแสดงถึงโหนดหรือจุดในกลุ่มคุณสมบัติซึ่งเขียนได้เป็น  $X = \{x[i] = 1, \dots, d\}$

นอกจากนี้ต้องมีการกำหนดการคำนวณค่ากึ่งกลางให้กับชุดคุณสมบัติแต่ละประเภทดังนี้  $\bar{C} = \{\bar{c}[i] : i = 1, \dots, d\}$  แทนจุดกึ่งกลางของกลุ่มคุณสมบัติ  $C$  โดยที่  $\bar{c}[i]$  แทนค่ากึ่งกลางของคุณสมบัติที่เป็นชนิดเดียวกันคือชนิดที่  $i$

ในการแบ่งกลุ่มของสมบัติที่เกี่ยวข้องกับการจราจรของแพ็คเกจเกิดในเครือข่ายจะต้องคำนึงถึงประเภทของคุณสมบัติที่กำลังพิจารณาโดยแบ่งประเภทของคุณสมบัติเป็น 3 ประเภทดังนี้ คุณสมบัตินเชิงเลขคณิต (Numerical Attribute), คุณสมบัตินเชิงประเภท (Categorical Attribute), คุณสมบัตินเชิงลำดับชั้น (Hierarchical Attribute)

#### 1) คุณสมบัตินเชิงเลขคณิต

คุณสมบัตินเชิงเลขคณิต คือ คุณสมบัติที่สามารถนำข้อมูลมาทำการคำนวณได้เลยโดยไม่ต้องผ่านการแทนค่าหรือเปลี่ยนแปลงค่าก่อน ตัวอย่างคุณสมบัตินที่ถือเป็นคุณสมบัตินเชิงเลขคณิตคือ จำนวนไบต์ข้อมูล คุณสมบัตินเชิงเลขคณิตสามารถเขียนได้เป็นปริมาณสเกลาร์  $x[i] \in R$  โดยมี  $\bar{c}[i]$  แทนค่ากึ่งกลางของคุณสมบัติที่เป็นคุณสมบัตินเชิงเลขคณิตของกลุ่มคุณสมบัติ  $C$  ที่มีจำนวนสมาชิกข้อมูล  $N$  ชุดคือ

$$\bar{c}[i] = \frac{1}{N} \sum_{j=1}^N x_j[i] \quad (1)$$

การคำนวณระยะห่างเชิงเลขคณิตคือการคำนวณระยะห่างระหว่างค่ากึ่งกลางข้อมูลของที่เป็นข้อมูลชนิดเลขคณิตเหมือนกันโดยอาศัยหลักการของยูคลิดคือการใช้ระยะห่างเมตริก  $\|\bullet\|$ :

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษานั้น ไม่ออบเลงไปใช้ประโยชน์ด้านการค้า  
 $d_n^i(C_1, C_2) = \|c_1[i] - c_2[i]\| = \left[ (c_1[i] - c_2[i])^2 \right]^{1/2}$  (2)  
 ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## 2) คุณสมบัติเชิงประเภท

คุณสมบัติเชิงประเภทคือคุณสมบัติที่ทำหน้าที่บ่งบอกชนิดหรือประเภทของแพ็คเก็ต ตัวอย่างคุณสมบัติเชิงประเภทเช่น โพรโทคอลและแพล็กของแพ็คเก็ต ในกรณีของคุณสมบัติ ชนิดนี้จะถูกแทนค่าด้วย  $x[i]$  ซึ่งเป็นเวกเตอร์ที่เป็นสมาชิกของจำนวนเต็ม  $Z^c$  โดยที่  $c$  คือ จำนวนค่าที่เป็นไปได้ของคุณสมบัตินั้น สำหรับเวกเตอร์  $x$  ซึ่งเป็นเวกเตอร์ของคุณสมบัติที่มี จำนวนค่าที่เป็นไปได้  $c$  ค่า จะสามารถเขียนแทนได้คือ  $x[i] = \{a_1, a_2, a_3, \dots, a_c\}$  โดยที่  $a_c \in \{0,1\}$  ซึ่งจะทำให้ผลรวมของ  $a_k$  ในหนึ่งเรคคอร์ดมีค่าเป็น 1 เสมอตามสมการ

$$\sum_{k=1}^c a_k = 1 \text{ เมื่อ } a_k \in x[i] \quad (3)$$

ค่ากึ่งกลางข้อมูลของคุณสมบัติเชิงประเภท  $\bar{c}[i]$  สามารถหาได้จากสมการ

$$\bar{c}[i] = \frac{1}{N} \sum_{j=1}^N x_j[i] = \left\{ \frac{A_1}{N}, \frac{A_2}{N}, \dots, \frac{A_c}{N} \right\} \quad (4)$$

เมื่อ

$$A_k = \sum_{j=1}^N x_{j,i} = \{a_{1,j}, a_{2,j}, \dots, a_{c,j}\}$$

ระยะห่างระหว่างกลุ่มข้อมูล  $C_1$  กับ  $C_2$  สำหรับคุณสมบัตินี้สามารถหาได้จากสมการ

$$d'_c(C_1, C_2) = \|\bar{c}_1 - \bar{c}_2\| = \frac{1}{N} \left[ \sum_{k=1}^c (A_{k,1} - A_{k,2})^2 \right]^{\frac{1}{2}} \quad (5)$$

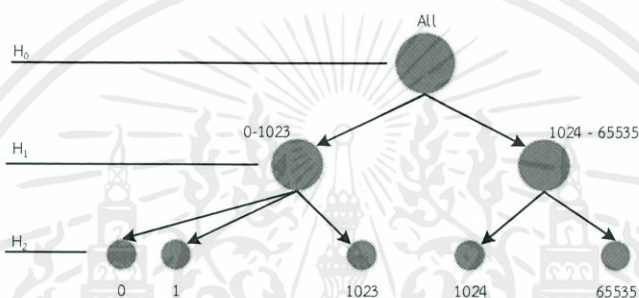
## 3) คุณสมบัติเชิงลำดับชั้น

ตัวอย่างคุณสมบัติของแพ็คเก็ตที่สามารถนำมาวิเคราะห์เป็นคุณสมบัติเชิงลำดับชั้น ได้แก่ หมายเลขพอร์ตและหมายเลขไอพี การนำคุณสมบัติเชิงลำดับชั้นมาคำนวณหรือวัด ความเหมือนกันนั้นจะพิจารณาคุณสมบัติเชิงลำดับชั้นในรูปแบบของต้นไม้ที่  $L$  ระดับนั้น เอกสารนี้เป็น ความหมายว่าต้นไม้จะมีความลึก  $L-1$  โดยที่โหนดที่อยู่ในระดับ  $l$  จะแทนด้วย  $H_l$  ดังนั้น ไม่ว่ากรณีใดๆ โหนดใดๆที่ไม่ได้ถือว่าเป็นลิฟโหนดซึ่งมีตำแหน่งอยู่ในระดับที่  $l$  จะต้องเป็นสมาชิกของ  $H_l$  ดังใช้ สมการ  $h_{l,i} \in H_l$  สำหรับต้นหมายที่มี  $L$  ระดับ การพิจารณาถึงค่าๆหนึ่งซึ่งเป็นค่าที่บอกถึง

การที่ปมสองปมใด ๆ มีบรรพบุรุษเดียวกันจะเรียกค่านั้นว่า Lowest Common Ancestor หรือ LCA ระหว่าง  $n_1$  กับ  $n_2$  แทนด้วย  $LCA(n_1, n_2)$  ส่วนความลึกของโหนด  $n$  ใดๆจะถูกนิยามว่าเป็นความยาวของเส้นทางตั้งแต่รากจนถึงโหนด  $n$

กำหนดให้  $path(n_1, n_2)$  แสดงถึงเส้นทางที่สั้นที่สุดระหว่างโหนด  $n_1$  ไปยัง โหนด  $n_2$  โหนดราก แสดงด้วย  $r$ , ลีฟโหนดแสดงด้วย  $l$  จะได้ว่า

$|path(l_1, l_2)| = |path(l_1, l_r)| + |path(l_r, l_2)|$  โดยที่  $|path(n_1, n_2)|$  หมายถึงความยาวเส้นทางระหว่างโหนด  $n_1$  ไปยังโหนด  $n_2$



รูปที่ 2.13 ตัวอย่างต้นไม้ทวิภาคสำหรับหมายเลขพอร์ต

จากรูปแสดงการนำหมายเลขพอร์ตแสดงในรูปแบบของต้นไม้ หมายเลขพอร์ตสามารถแสดงในรูปแบบต้นไม้โดยมีเพียง 2 ระดับ ( $L=2$ ) อยู่บนโดเมนของเลขจำนวนเต็ม จากรูปจะเห็นได้ว่าในระดับที่ 0 นั้น  $H_0$  แสดงถึงทุกๆหมายเลขพอร์ต (ตั้งแต่  $0 - 2^{16}$ )  $H_0 = \{All\}$  หมายความว่าทุกค่าหมายเลขพอร์ตถือเป็นกลุ่มเดียวกัน ส่วนในระดับที่ 1 นั้นจะมีค่าสองค่าที่บอกความแตกต่างของโหนดลูกทั้งสองคือ “Low” และ “High” ซึ่งก็คือ  $H_1 = \{“Low”, “High”\}$  นั้นหมายความว่าถ้าหมายเลขพอร์ตเป็น 1 ซึ่งจะถูกจัดอยู่ในกลุ่ม “Low” ส่วนในกรณีของระดับที่ 2 นั้นจะได้ว่า  $H_2 = \{0,1,2,\dots,65535\}$  หมายความว่า จะมี 65536 กลุ่มแต่ละกลุ่มมีสมาชิกก็คือตัวเองนั่นเอง ในกรณีของหมายเลขไอพีก็เช่นเดียวกันคือสามารถเขียนให้อยู่ในรูปแบบต้นไม้ได้

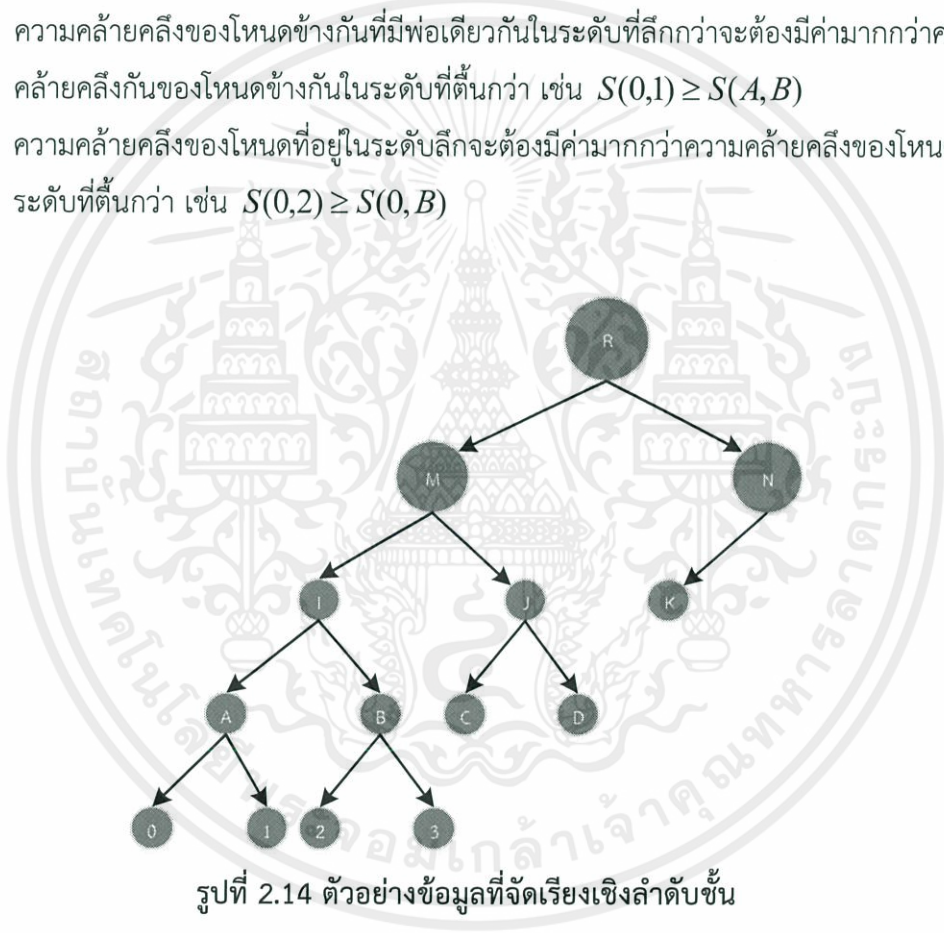
### 2.7.1 คุณสมบัติที่จำเป็นสำหรับการวัดความคล้ายคลึงระหว่างข้อมูลแบบลำดับชั้น

การนำข้อมูลมาจัดเรียงในรูปแบบต้นไม้ต้องทราบคุณสมบัติที่จำเป็นสำหรับการวัดความคล้ายคลึงระหว่างข้อมูลแบบลำดับชั้น หรือ Desired Properties of a Hierarchical Similarity Measure มิฉะนั้นจะไม่สามารถใช้กระบวนการทางคณิตศาสตร์แบบปกติวัดความคล้ายกันของแต่ละโหนดได้ ดังนั้นต้องมีการกำหนดกระบวนการในการวัดความคล้ายกันของแต่ละโหนดนั้น กำหนดให้  $S$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งยังมีเหตุผลเบื้องเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีกรนำไปใช้

คือกระบวนการวัดความคล้ายคลึงระหว่างข้อมูลแบบลำดับชั้น คุณสมบัติที่จำเป็นของกระบวนการดังกล่าวมีดังนี้

- 1) S ควรจะเป็นกระบวนการที่สมมาตรสำหรับสองโหนด I และ J ใดๆ  $S(I, J) = S(J, I)$
- 2) ความคล้ายคลึงระหว่างโหนด X กับโหนดพ่อ Y จะต้องมีค่ามากกว่าโหนด X กับโหนดข้างเคียงโหนด Y เช่น  $S(0, A) > S(0, B)$
- 3) ค่าของ S จะต้องมีการลดลงเมื่อมีการวัดความคล้ายคลึงในแนวตั้งเริ่มจากโหนด X กับตัวเองขึ้นไปเรื่อยๆจนถึงราก เช่น  $S(0,0) > S(0, A) > S(0, I) > S(0, M)$
- 4) ความคล้ายคลึงของโหนดข้างกันที่มีพ่อเดียวกันในระดับที่ลึกกว่าจะต้องมีค่ามากกว่าความคล้ายคลึงกันของโหนดข้างกันในระดับที่ตื้นกว่า เช่น  $S(0,1) \geq S(A, B)$
- 5) ความคล้ายคลึงของโหนดที่อยู่ในระดับลึกจะต้องมีค่ามากกว่าความคล้ายคลึงของโหนดในระดับที่ตื้นกว่า เช่น  $S(0,2) \geq S(0, B)$



รูปที่ 2.14 ตัวอย่างข้อมูลที่จัดเรียงเชิงลำดับชั้น

ในการวัดความคล้ายคลึงระหว่างข้อมูลแบบลำดับชั้นนั้นจะต้องกำหนดการดำเนินการทางคณิตศาสตร์ที่รองรับคุณสมบัติที่กล่าวมาทั้งห้าข้อ จากการศึกษาพบว่า Hierarchical Similarity Measure 1 (HSM<sub>1</sub>) ตามเอกสารอ้างอิง<sup>[1]</sup> มีคุณสมบัติครบถ้วนซึ่งเขียนเป็นสมการได้คือ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้า S<sub>HSM<sub>1</sub></sub>(n<sub>1</sub>, n<sub>2</sub>) =  $\frac{path(root, n_1) + path(root, n_2)}{L}$  ไม่อนุญาตให้นำไปใช้

## 2.7.2 การวิเคราะห์หมายเลขไอพีแอดเดรสแบบสมบัติเชิงลำดับชั้น

การวิเคราะห์หมายเลขไอพีสามารถทำได้โดยการจัดการกับหมายเลขไอพีด้วยต้นไม้แบบทวิภาค (Binary Tree) ของพรีฟิกส์ ซึ่งมี 32 ระดับตามรูปแบบของหมายเลขไอพีที่เป็นเลขฐานสองขนาด 32 บิตซึ่งทำให้หมายเลขไอพีมีทั้งสิ้น  $2^{32}$  หมายเลข ค่าพรีฟิกส์ของหมายเลขไอพีจะแสดงถึงระดับในต้นไม้ทวิภาค

**คำนิยาม 1** พรีฟิกส์ของหมายเลขไอพี (Prefix of an IP address)

กำหนดหมายเลขไอพีซึ่งเป็นเลขฐานสองขนาด 32 บิตและพรีฟิกส์เป็นเลขฐานสองขนาด  $p$  บิตจะได้ว่า พรีฟิกส์คือ  $p$  บิตแรกที่มีนัยสำคัญสูงสุดของหมายเลขไอพีเขียนแทนด้วย  $IP/p$  เมื่อ  $0 \leq p \leq 32$

**คำนิยาม 2** CommonPrefix and AggregateIP

CommonPrefix ถูกนิยามว่าเป็นพรีฟิกส์  $p$  ที่ยาวที่สุดที่ทำให้  $IP_1/p = IP_2/p$  ซึ่งเปรียบเทียบกับได้กับสมการ  $|path(root, n_1) \cap path(root, n_2)|$

AggregateIP ของหมายเลขไอพีสองหมายเลข  $IP_1$  และ  $IP_2$  ใดๆถูกนิยามว่าเป็น  $IP/p$  เมื่อ  $p$  เป็น CommonPrefix ของ  $IP_1$  และ  $IP_2$  ตัวอย่างเช่น

$$AggregateIP(203.192.32.127, 203.192.32.128) = 203.192.32.0/24$$

จากนิยามทำให้สามารถกำหนดการวัดระยะห่างเชิงลำดับชั้นของหมายเลขไอพีได้ดังสมการ

$$D_{HSM_1}(n_1, n_2) = 1 - S_{HSM_1}(n_1, n_2)$$

$$D_{HSM_1}(n_1, n_2) = (L - |path(root, n_1) \cap path(root, n_2)|) / L$$

ในกลุ่มของ  $C$  สำหรับสมบัติเชิงลำดับชั้นนั้นสามารถถูกพิจารณาค่ากึ่งกลางของข้อมูลได้คือค่า  $\overline{IP}/p$  ซึ่งเป็นค่า AggregateIP ของกลุ่มหมายเลขไอพีในทางปฏิบัติแล้วจะไม่สามารถหาค่าซบเนตที่มีค่าน้อยกว่า 8 ได้อยู่แล้วดังนั้นหมายเลขไอพีสองหมายเลขใดๆที่มี CommonPrefix  $p$  ที่น้อยกว่าหรือเท่ากับ 8 สามารถพิจารณาว่าไม่มีความคล้ายคลึงกันเลยได้และจะสามารถหาระยะห่างระหว่างข้อมูลสองกลุ่มระหว่าง  $C_1$  และ  $C_2$  ได้ดังนี้

$$d'_n(C_1, C_2) = \begin{cases} (32-p)/24 & p > 8 \\ 1 & p \leq 8 \end{cases}$$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

### 2.7.3 การรวมผลต่างระยะห่างเชิงเลขคณิต, เชิงประเภท และเชิงลำดับชั้น

กำหนดให้ระยะห่างของข้อมูลคุณสมบัติเชิงเลขคณิตคือ  $d_n$ , คุณสมบัติเชิงประเภท  $d_c$  และคุณสมบัติเชิงลำดับชั้นคือ  $d_h$  จะสามารถสรุประยะห่างระหว่างกลุ่มข้อมูล  $C_1$  และ  $C_2$  ได้คือ

$$D_i(C_1, C_2) = \begin{cases} d_n^i(C_1, C_2), i \in \text{numerical attribute}, \\ d_c^i(C_1, C_2), i \in \text{categorical attribute}, \\ d_h^i(C_1, C_2), i \in \text{hierarchical attribute} \end{cases} \quad (6)$$

และจะได้ระยะห่างระหว่างกลุ่มด้วยผลรวมของระยะห่างของทุกคุณสมบัติเป็น

$$D(C_1, C_2) = \left[ \sum_{i=1}^d D_i(C_1, C_2) \right]^{\frac{1}{2}} \quad (7)$$

### 2.7.4 การคำนวณรัศมีของกลุ่มข้อมูล

ในการควบคุมความแปรปรวนของข้อมูลในกลุ่มข้อมูลนั้นเราต้องมีการคำนวณหรือกำหนดขอบเขตบางอย่างของข้อมูลไม่ให้ข้อมูลภายในกลุ่มมีค่ามากเกินไป ในกรณีของสมบัติเชิงเลขคณิตและเชิงประเภทนั้นเราสามารถใช้ในการคำนวณทางคณิตศาสตร์ซึ่งจะได้กล่าวต่อไปในการกำหนดขอบเขตของข้อมูลได้ แต่สำหรับสมบัติเชิงลำดับชั้น เช่น หมายเลขไอพีหรือหมายเลขพอร์ตจะต้องใช้กระบวนการที่แตกต่างออกไป ผลรวมทางคณิตศาสตร์ของรัศมีของสมบัติแต่ละประเภทนั้นก็คือรัศมีของกลุ่มข้อมูลนั่นเอง ขั้นตอนการคำนวณรัศมีของกลุ่มข้อมูลมีดังต่อไปนี้

- 1) รัศมี  $R_n$  สำหรับสมบัติเชิงเลขคณิตที่  $i$  ในกลุ่มข้อมูลที่มี  $N$  เรคคอร์ดและค่ากึ่งกลางข้อมูล  $\bar{c}[i]$  คือ

$$R_n = \sqrt{\frac{1}{N} \sum_{j=1}^N (x_j[i])^2 - (\bar{c}[i])^2}$$

- 2) รัศมี  $R_c$  สำหรับคุณสมบัติเชิงประเภทที่  $i$  ซึ่งมีค่าที่สามารถเป็นไปได้  $c$  ค่าคือ

$$R_n = \sqrt{\frac{1}{N} \sum_{j=1}^N \left( \sum_{k=1}^c (a_{k,j})^2 \right) - \frac{1}{c^2} \left( \sum_{k=1}^c A_k \right)^2}$$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับบุคคลอื่นที่สนใจศึกษาและนำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- 3) สำหรับสมบัติเชิงลำดับชั้นนั้นไม่สามารถใช้การคำนวณเช่นเดียวกับข้อ ก. และ ข. เนื่องจากสมบัติประเภทนี้มีการจัดเรียงในรูปแบบเป็นลำดับชั้น การใช้ค่ามากที่สุดและน้อยที่สุดสำหรับพิจารณารัศมีจึงเป็นสิ่งที่เหมาะสม ในการพิจารณารัศมีของกลุ่มข้อมูลประเภทนี้จะใช้สมการ

$$R_i = (32 - \text{Common Prefix}(\min IP, \text{MaxIP})) / 32 \quad (9)$$

ตัวอย่างเช่น

$$C_1.\text{range} = (192.168.0.1, 192.168.0.2) \quad \text{Radius} = (32-31)/32 = 1/32$$

$$C_2.\text{range} = (192.168.0.1, 192.168.0.255) \quad \text{Radius} = (32-24)/32 = 1/4$$

## 2.8 เน็ตโฟลว์

เน็ตโฟลว์ (NetFlow) เป็นโพรโทคอลเครือข่ายที่ถูกพัฒนาโดยระบบของซิสโก้ (Cisco) ใช้ในการเก็บรวบรวมข้อมูลจากการทราฟฟิกไอพี เน็ตโฟลว์นั้นถูกใช้โดยผู้เชี่ยวชาญด้านไอทีในการวิเคราะห์การไหลของทราฟฟิกและกำหนดปริมาณทราฟฟิกที่กำลังมาว่าจะไปที่ไหนและมีจำนวนเท่าไรที่ถูกสร้างขึ้น

เน็ตโฟลว์จะเปิดใช้งานเราเตอร์ในการนำข้อมูลทราฟฟิกและสถิติออกมาบันทึกไว้ซึ่งจะถูกเก็บรวบรวมโดยตัวเก็บรวบรวมเน็ตโฟลว์หรือเน็ตโฟลว์คอลเล็คเตอร์ (NetFlow Collector) ซึ่งตัวเก็บจะมีการวิเคราะห์ทราฟฟิกจริงและนำเสนอให้กับผู้ใช้ เน็ตโฟลว์ได้กลายเป็นมาตรฐานอุตสาหกรรมสำหรับตรวจสอบทราฟฟิกและได้รับการสนับสนุนจากหลายแพลตฟอร์ม การไหลของข้อมูลในเครือข่าวนั้นสามารถกำหนดได้หลายทาง ในซิสโก้มาตรฐานเน็ตโฟลว์เวอร์ชันที่ 5 มีการกำหนดการไหลลำดับของแพ็คเก็ตหลายค่า เช่น

- ไอพีแอดเดรสต้นทาง, ไอพีแอดเดรสปลายทาง
- ไอพีแอดเดรสของ next hop เราท์เตอร์
- ขนาดของแพ็คเก็ต
- ทีซีพี / ยูดีพี ของพอร์ตต้นทาง และพอร์ตปลายทาง
- ประเภทการให้บริการของไอพี
- ชนิดของโพรโทคอลไอพี

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## 2.9 เครื่องมือที่ใช้ในการดำเนินงาน

เครื่องมือที่ใช้ในการดำเนินงานมีดังนี้

### 1) เน็ตโฟลว์ อนาไลเซอร์ (NetFlow Analyzer)

เน็ตโฟลว์ อนาไลเซอร์ เป็นอุปกรณ์ที่ใช้ในการวิเคราะห์ทราฟฟิก และรายงานเกี่ยวกับที่เครือข่ายที่มีการใช้งาน เพื่อแสดงผลตามเวลาจริงในระบบเครือข่ายอย่างมีประสิทธิภาพ คุณสมบัติที่สำคัญ เช่น

- ตรวจสอบแอปพลิเคชันและโพรโทคอล ดูการใช้งานบนเครือข่ายที่มีการใช้งานบ่อยๆ
- รายงานทราฟฟิกตามเวลาจริง รับการแสดงผลเวลาจริงของทราฟฟิกในเครือข่าย
- การเตือนสามารถตั้งค่าการเตือนตามค่าเทรชโฮลสำหรับทราฟฟิกบนเครือข่าย

### 2) เอฟโพรบ (FProbe)

เป็นเครื่องมือลิบแคปที่ใช้ในการเก็บรวบรวมข้อมูลทราฟฟิกของเครือข่ายและปล่อยทราฟฟิคนั้นออกมาในรูปแบบเน็ตโฟลว์ไปยังที่เก็บข้อมูลที่ระบุไว้ เอฟโพรบรองรับการทำงานของ Netflow เวอร์ชันที่ 5 ซึ่งเป็นเวอร์ชันมาตรฐานสำหรับอุปกรณ์เครือข่ายทุกตัวที่สามารถทำงานกับ Netflow ได้ เอฟโพรบทำงานบนระบบปฏิบัติการลินุกซ์ ซึ่งจะทำการรวบรวมข้อมูลจากการใช้งานเครือข่ายที่ผ่านเครื่องคอมพิวเตอร์ที่ได้ติดตั้งเอฟโพรบไว้และจะทำการจัดทำข้อมูลสรุปเป็นแพ็คเกจส่งไปยังเครื่องปลายทางที่ได้ตั้งค่าไว้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## บทที่ 3

# การออกแบบและการพัฒนา

ในบทนี้จะกล่าวถึงการนำความรู้จากการค้นคว้ามาออกแบบระบบตรวจหาบอทเน็ตจากการวิเคราะห์พฤติกรรมของระบบเครือข่าย โดยระบบจะทำหน้าที่ในการบันทึกการรับส่งข้อมูลต่างๆบนเครือข่ายที่ซีพี/ไอพีลงในฐานข้อมูล จากนั้นจะนำข้อมูลที่เก็บไว้มาแบ่งกลุ่มตามลักษณะของการรับส่งข้อมูลเพื่อเปรียบเทียบกับพฤติกรรมของบอทเน็ตในภายหลัง

### 3.1 รายละเอียดของระบบที่พัฒนา

การออกแบบและการพัฒนาระบบ จะต้องพิจารณาถึงรายละเอียดของระบบ ขั้นตอนการทำงานต่างๆ ของโปรแกรม รวมไปถึงการแสดงผลลัพธ์ของการวิเคราะห์ของโปรแกรม

#### 3.1.1 รายละเอียดการนำเข้าข้อมูล (Input Specification)

การรับ-ส่งข้อมูลระหว่างเครื่องคอมพิวเตอร์ในเครือข่ายกับอินเทอร์เน็ต เป็นข้อมูลนำเข้าพื้นฐานของระบบซึ่งจะอาศัยความสามารถของ เน็ตโพล์ Generator สร้างแพ็คเกจเน็ตโพล์ที่สรุปการส่งข้อมูลระหว่างแต่ละเครื่องคอมพิวเตอร์เพื่อนำข้อมูลที่ได้ไปทำการจัดกลุ่มข้อมูลแล้วเปรียบเทียบความใกล้เคียงกันกับข้อมูลของบอทเน็ต

#### 3.1.2 รายละเอียดผลลัพธ์ของระบบ (Output Specification)

ผลลัพธ์ของระบบสามารถบ่งบอกได้ว่าการใช้งานของเครื่องคอมพิวเตอร์ใดในช่วงเวลาใด โดยผลลัพธ์ที่ได้อาจจะเป็นหรือไม่เป็นบอทเน็ตได้ เพราะการวิเคราะห์การใช้งานเครือข่ายอาศัยการแบ่งกลุ่มพฤติกรรมของการใช้งานเครือข่ายซึ่งต่างจากการตรวจหาบอทเน็ตจากลักษณะเฉพาะ (Signature) ซึ่งจะสามารถบ่งบอกได้เลยว่าเป็นบอทเน็ตหรือไม่ แต่ข้อดีของการวิเคราะห์พฤติกรรมคือสามารถประเมินได้ว่าเครื่องใดอาจเป็นบอทเน็ตโดยไม่ต้องอาศัยการอัปเดตฐานข้อมูลลักษณะเฉพาะของบอทเน็ต

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

### 3.1.3 ขอบเขตของระบบที่พัฒนา

- 1) ระบบตรวจจับบอทเน็ตที่พัฒนาขึ้นเพื่อเป็นต้นแบบในการตรวจจับบอทเน็ตเชิงพฤติกรรมซึ่งจะต้องอาศัยการคำนวณและแบ่งกลุ่มของพฤติกรรมการทำงานของข้อมูลปริมาณมากซึ่งจำเป็นต้องใช้ความสามารถของเครื่องเซิร์ฟเวอร์ที่สูง
- 2) ระบบนี้ได้ถูกพัฒนามาโดยใช้การวิเคราะห์ข้อมูลจากการสังเกตทำให้ไม่สามารถตรวจจับบอทเน็ตแบบทันที (Real Time) ได้
- 3) ระบบถูกพัฒนาขึ้นโดยใช้ภาษาจาวาซึ่งมีความเร็วน้อยกว่า Native Language เช่น C หรือ C++

### 3.1.4 ข้อจำกัดของระบบที่พัฒนา

- 1) ระบบนี้เป็นระบบที่สามารถตรวจจับบอทเน็ตได้จากพฤติกรรมที่ใกล้เคียงตัวอย่างบอทเน็ตที่ได้นับทีวกไว้ซึ่งจะไม่สามารถตรวจจับบอทเน็ตที่มีพฤติกรรมที่แตกต่างจากที่ระบุไว้มากๆได้
- 2) ระบบนี้ไม่สามารถตรวจจับบอทเน็ตแบบ Real Time ได้
- 3) ไม่สามารถตรวจจับบอทเน็ตที่ไม่แสดงพฤติกรรมออกมาได้

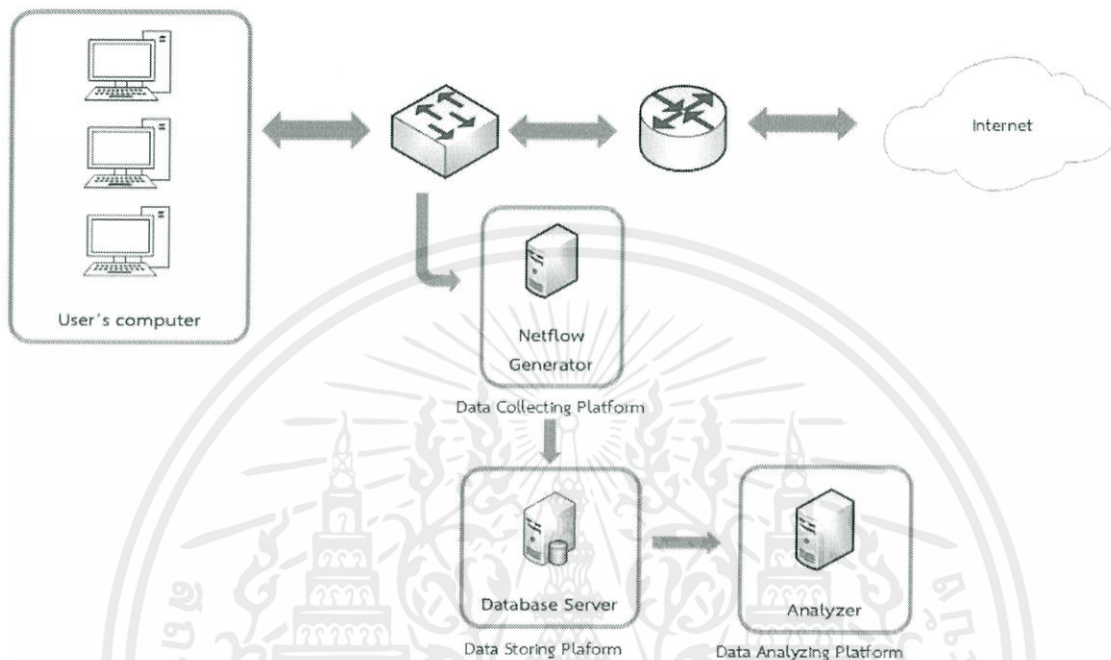
### 3.1.5 เครื่องมือที่ใช้ในการพัฒนา

- 1) สภาพแวดล้อมในการพัฒนา
  - หน่วยประมวลผล Intel® Core™ i5-2410M 2.30 GHz
  - หน่วยความจำหลัก 8 GB
  - หน่วยความจำรอง 640 GB
  - ระบบปฏิบัติการ Windows
- 2) ซอฟต์แวร์ที่ใช้พัฒนา
  - Eclipse Standard 4.3
  - mysql-connector-java-5.1.26-bin.jar
  - Fprobe
- 3) ภาษาที่ใช้ในการพัฒนา คือ ภาษา Java

## 3.2 รูปแบบโครงสร้างของระบบ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า โครงสร้างของระบบประกอบไปด้วยโครงสร้างหลักสามส่วนคือส่วนรวบรวมข้อมูลทำหน้าที่ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งยังมีชุดคำสั่งและต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งหากนำไปใช้ในการเก็บรวบรวมข้อมูลการรับ-ส่งข้อมูลของระบบเครือข่าย ส่วนเก็บรักษาข้อมูลทำหน้าที่ในการ

เก็บรักษาข้อมูลที่ได้จากการรวบรวมข้อมูลไว้ในฐานข้อมูล และส่วนวิเคราะห์ข้อมูลและแสดงผลทำหน้าที้นำข้อมูลจากฐานข้อมูลมาวิเคราะห์เปรียบเทียบกับข้อมูลของบอทเน็ตตัวอย่าง



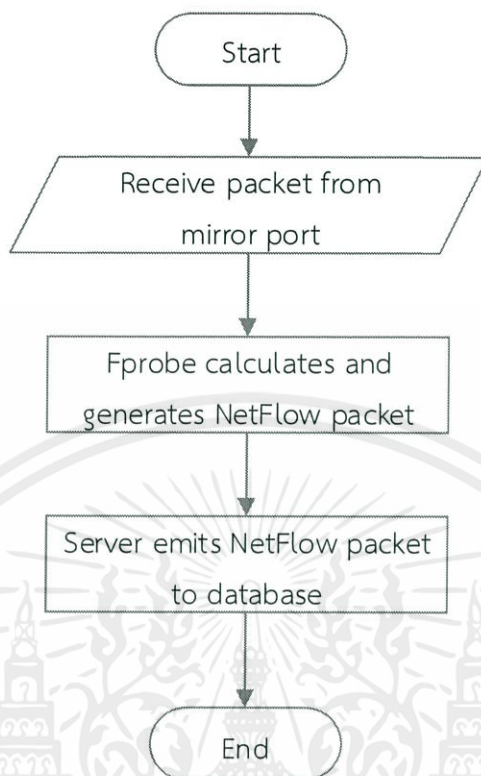
รูปที่ 3.1 ภาพรวมโครงสร้างของระบบ

### 3.2.1 ส่วนรวบรวมข้อมูล

ส่วนรวบรวมข้อมูล (Data Collecting Platform) สามารถใช้ฟังก์ชันพิเศษที่มีในสวิตช์หรือเราท์เตอร์ของซิสโก้หรือของผู้ผลิตรายอื่นๆที่รองรับการใช้งานเน็ตโพล์ได้ แต่เนื่องจากอุปกรณ์เครือข่ายของสาขาวิชาไม่รองรับการใช้งานฟังก์ชันดังกล่าวจึงจำเป็นต้องหาวิธีอื่นในการสร้างแพ็คเก็ตเน็ตโพล์ขึ้นมา

ส่วนรวบรวมข้อมูลประกอบด้วยเครื่องเซิร์ฟเวอร์ที่ได้ทำการลงซอฟต์แวร์เอพโพรซึ่งเป็นซอฟต์แวร์ที่สามารถสร้างแพ็คเก็ตเน็ตโพล์ได้จากการวิเคราะห์แพ็คเก็ตต่างๆที่วิ่งผ่านเครื่องเซิร์ฟเวอร์ ตัวเซิร์ฟเวอร์ทำหน้าที่รวบรวมข้อมูลจากอุปกรณ์เครือข่ายสวิตช์หรือเราท์เตอร์ผ่านมัลเลอร์พอร์ตเพื่อนำข้อมูลที่ได้ไปสรุปและสร้างเป็นแพ็คเก็ตเน็ตโพล์และส่งต่อไปยังส่วนเก็บรักษาข้อมูล

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



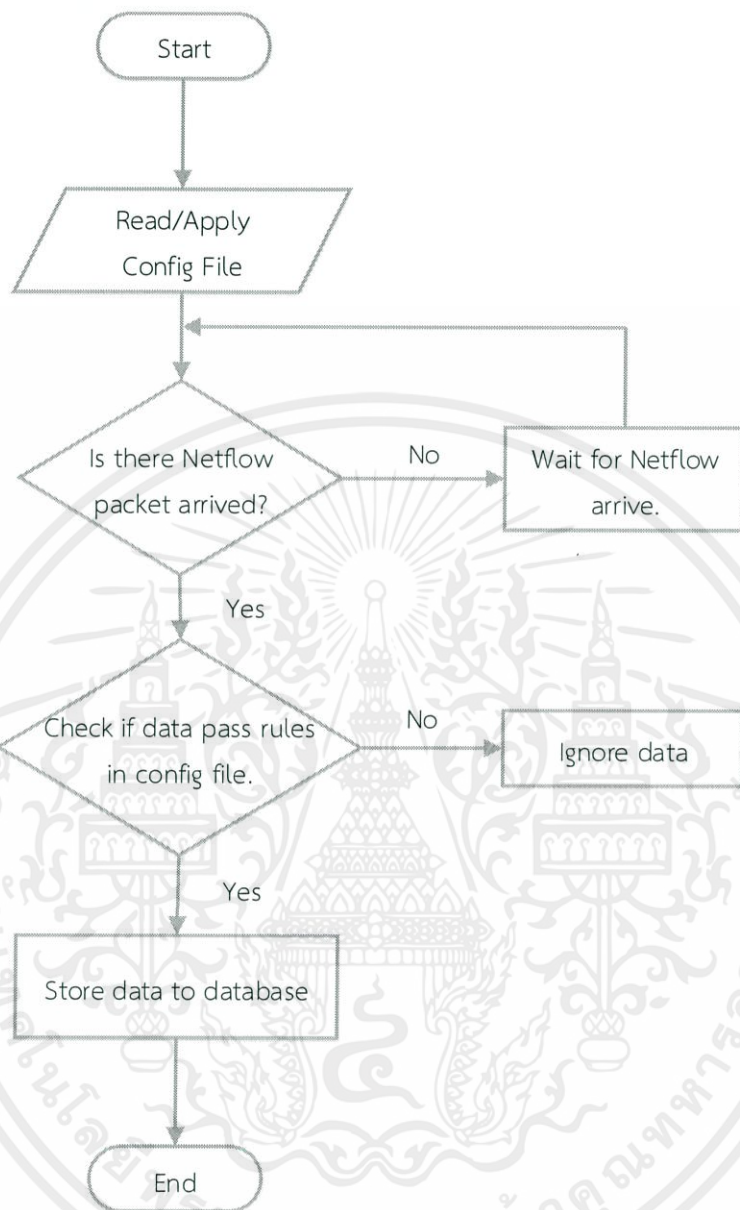
รูปที่ 3.2 ผังแสดงการทำงานของส่วนรวบรวมข้อมูล

### 3.2.2 ส่วนเก็บรักษาข้อมูล

หลังจากที่ส่วนรวบรวมข้อมูลได้ทำการรวบรวมข้อมูลการใช้งานเครือข่ายแล้วก็สร้างแพ็คเกจเน็ตฟลิวส่งมายังส่วนเก็บรักษาข้อมูล

ส่วนเก็บรักษาข้อมูล (Data Storing Platform) คือส่วนที่ทำหน้าที่ในการรักษาข้อมูลสองประเภท คือ ข้อมูลดิบที่ได้จากส่วนรวบรวมข้อมูลและข้อมูลที่ผ่านการคำนวณแล้ว ส่วนเก็บรักษาข้อมูลจะใช้ซอฟต์แวร์ MySQL ซึ่งเป็น รีเลชันเนล ดาตาเบส (Relational Database) ชนิดหนึ่งที่ได้รับ ความนิยมและมีความรวดเร็วในการเก็บและดึงข้อมูลมาใช้งานทำหน้าที่เก็บรักษาข้อมูลซึ่งเหตุผลที่เลือกใช้ฐานข้อมูลนี้ได้กล่าวไว้แล้วในบทที่สอง การที่เลือกใช้ รีเลชันเนล ดาตาเบส สาเหตุที่สำคัญอย่างหนึ่งคือการทำงานที่ต้องเก็บข้อมูลขนาดเล็กเป็นปริมาณมากซึ่งจำเป็นต้องใช้ความสามารถของ ดีบีเอ็มเอส ที่รวดเร็วในการดึงข้อมูลที่ต้องการมาใช้ในการคำนวณ

ในการเก็บรักษาข้อมูลที่ได้จากส่วนรวบรวมข้อมูล แทนที่จะเก็บข้อมูลทั้งหมดที่ได้จากส่วนรวบรวมข้อมูลระบบจะมีการกรองแพ็คเกจ (Data Reduction Filters) ที่ไม่จำเป็นในการนำไปวิเคราะห์หรือออกไปก่อนเช่น แพ็คเกจที่มีการระบุที่อยู่ปลายทางเป็นที่อยู่ภายในเน็ตเวิร์คเดียวกันหรือเป็นที่อยู่ที่อยู่จึกที่อยู่แล้ว เช่น กูเกิ้ล, ยาฮู, ยูทูบ และอื่นๆ

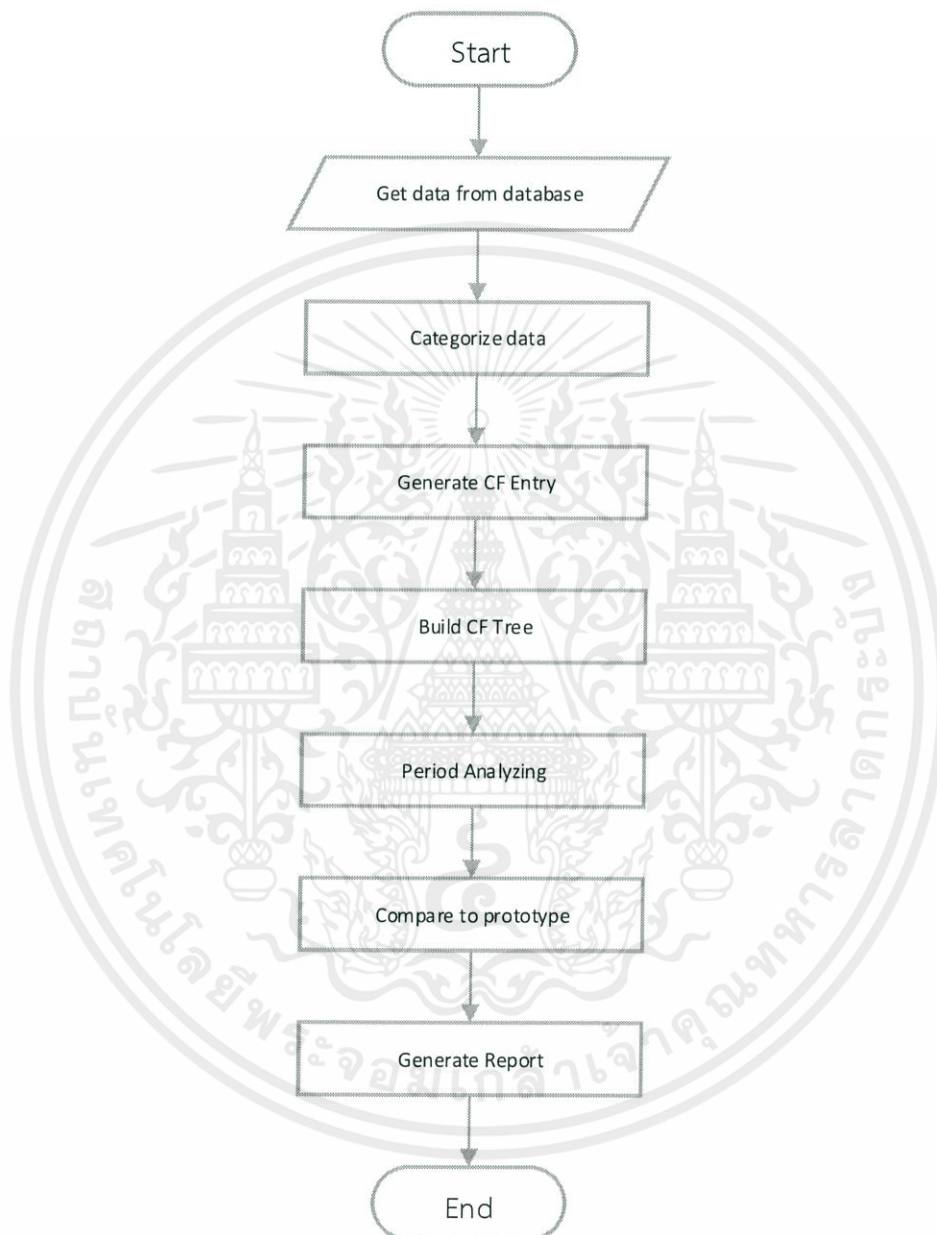


รูปที่ 3.3 ผังแสดงการทำงานของส่วนเก็บรักษาข้อมูล

### 3.2.3 ส่วนวิเคราะห์ข้อมูล

ส่วนวิเคราะห์ข้อมูล (Data Analyzing Plaform) จะทำงานภายหลังจากที่ได้ข้อมูลมาเก็บลงในฐานข้อมูลแล้วก็จะสามารถนำข้อมูลมาวิเคราะห์ได้ ซึ่งข้อมูลที่น่ามาวิเคราะห์นั้นเป็นข้อมูลดิบที่ได้เก็บข้อมูลจากส่วนรวบรวมข้อมูล ปัญหาอย่างหนึ่งของการนำข้อมูลเกี่ยวกับการใช้งานเครือข่ายมาวิเคราะห์หรือคำนวณในทางคณิตศาสตร์คือข้อมูลที่ได้มานั้นส่วนใหญ่ไม่ได้เป็นข้อมูลเชิงคำนวณได้ (Numerical) แต่ถึงอย่างไรก็ตามก็ยังสามารถนำข้อมูลดังกล่าวมาคำนวณได้โดยผ่านกระบวนการใช้

แทนค่าข้อมูลด้วยชุดข้อมูลที่เหมาะสมและต้องอาศัยกระบวนการในการคำนวณที่แตกต่างไปจากการคำนวณแบบปกติ กระบวนการดังกล่าวได้ถูกกล่าวถึงแล้วในบทที่ 2



รูปที่ 3.4 ผังแสดงการทำงานของส่วนวิเคราะห์ข้อมูล

### 3.3 แนวคิดการวิเคราะห์ข้อมูล

จากกระบวนการทางคณิตศาสตร์ในบทที่ 2 ในหัวข้อนี้จะเป็นการนำเอาวิธีการต่างๆ มาใช้ในการวิเคราะห์การจราจรในเครือข่าย โดยใช้พื้นฐานของการแบ่งกลุ่มข้อมูลเชิงลำดับชั้น เพื่อสรุป

พฤติกรรมของแพ็คเก็ตส่วนใหญ่ที่สามารถสังเกตได้จากการใช้งานเครือข่าย แนวคิดของโครงการจะนำเอาข้อมูลของแต่ละแพ็คเก็ตมาวิเคราะห์และสร้างเป็นต้นไม้เชิงลำดับชั้นตามอัลกอริทึมเบิร์ชโดยแต่ละเรคคอร์ดของข้อมูลที่ได้จากเครือข่ายจะถูกเพิ่มเข้าไปในกลุ่มข้อมูลที่มีความใกล้เคียงกับเรคคอร์ดที่พิจารณาที่สุดโดยใช้การวัดผลต่างเชิงเลขคณิต, เชิงประเภทและเชิงลำดับชั้นสำหรับทุกๆคุณสมบัติ โดยที่แต่ละเรคคอร์ดจะถูกรวมเข้าไปในกลุ่มข้อมูลหรือไม่ขึ้นอยู่กับสองปัจจัยคือ รัศมีของกลุ่มข้อมูล ( $R$ ) ซึ่งจะอธิบายวิธีการรวมในภายหลังและจำนวนข้อมูลมากที่สุดที่สามารถอยู่ในกลุ่มข้อมูลนั้นๆ ได้

### 3.3.1 รูปแบบการแบ่งกลุ่มข้อมูล

โครงการนี้จะทำการแบ่งกลุ่มข้อมูลตามแนวทางของอัลกอริทึมเบิร์ชซึ่งจะทำให้ได้ผลลัพธ์เป็นต้นไม้เชิงลำดับชั้นหรือซีเอฟทีรีลีฟโหนดของซีเอฟทีรีลีฟจะแสดงถึงรายละเอียดของข้อมูลที่แบ่งย่อยที่สุดในแต่ละกลุ่มข้อมูลและรากของ ซีเอฟทีรีลีฟ จะแสดงถึงคุณลักษณะทั่วไปของกลุ่มข้อมูลนั้นๆ กลุ่มข้อมูลใดๆ  $C_i$  จะถูกแสดงด้วยเวกเตอร์ ซีเอฟทีรีลีฟ ซึ่งมีข้อมูลสรุปที่เพียงพอที่จะใช้คำนวณจุดกึ่งกลางของข้อมูล  $\bar{c}_i$  และรัศมี  $R_i$  ของกลุ่มข้อมูล ในทำนองเดียวกับกรณีของโหนดพ็อก็สามารถใช้ข้อมูลของโหนดหลายๆในการคำนวณจุดกึ่งกลางข้อมูลและรัศมีของข้อมูลในกลุ่มข้อมูลเดียวกันได้ ขั้นตอนในการสร้างซีเอฟทีรีลีฟคือเมื่อมีข้อมูลเรคคอร์ดใหม่  $X$  จากเครือข่าย เรคคอร์ด  $X$  จะถูกดำเนินการเพิ่มเข้าไปในกลุ่มข้อมูล  $C_i$  ที่ใกล้เคียงกับเรคคอร์ด  $X$  ที่สุดโดยใช้การวัดระยะห่าง  $D(X, C_i)$  ในสมการที่ (7) ถ้าการเพิ่ม  $X$  เข้าไปในกลุ่มข้อมูล  $C_i$  แล้วทำให้รัศมีของกลุ่มข้อมูลที่คำนวณใหม่มีค่าเกินกว่าค่าเทรชโฮลด์  $T$  ซึ่งมีค่าอยู่ในช่วง  $[0, 1]$  จะต้องมีการสร้างกลุ่มข้อมูลใหม่ซึ่งอยู่บนพื้นฐานของข้อมูลเรคคอร์ด  $X$  แต่ถ้ารัศมีที่คำนวณได้ใหม่มีค่าอยู่ไม่เกินเทรชโฮลด์  $T$  ก็จะต้องมีการตรวจสอบว่าจำนวนข้อมูลในกลุ่มข้อมูลมีค่ามากที่สุด  $m$  หรือไม่ ถ้าไม่ก็สามารถเพิ่มเรคคอร์ด  $X$  เข้าไปในกลุ่มข้อมูลได้ ในกรณีที่ต้องมีการสร้างกลุ่มข้อมูลใหม่สำหรับเรคคอร์ด  $X$  แต่โหนดที่รองรับกลุ่มข้อมูลนั้นมีจำนวนกลุ่มข้อมูลเต็มที่แล้วก็ต้องมีการสร้างโหนดใหม่ขึ้นมาเพื่อรองรับเรคคอร์ด  $X$  แต่อย่างไรก็ตามเรคคอร์ด  $X$  จะไม่ถูกจัดให้อยู่ในกลุ่มข้อมูลใหม่ในโหนดใหม่ทันทีแต่ต้องมีการจัดเรียงกลุ่มข้อมูลใหม่โดยแบ่งแยกกลุ่มข้อมูลที่มีระยะห่างข้อมูลมากที่สุดออกมาอยู่คนละกลุ่มก่อนจากนั้นจะมีการเพิ่มเวกเตอร์ซีเอฟทีรีลีฟเข้าไปยังโหนดพ็อสำหรับโหนดใหม่ที่ถูกสร้างขึ้นมา

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

1. find the nearest cluster  $C_i$  to record  $X$  in  $N$ 
      2. if  $N$  is a leaf node then
3.      $C_i' \leftarrow$  insert  $X$  into  $C_i$ 
4.     Calculate radius  $R_i$  of  $C_i'$ 
5.     if  $R_i \leq T$ , where  $T$  is a threshold in  $[0,1]$  then
6.         update  $C_i \leftarrow C_i'$ 
7.     else  $\{R > T\}$ 
8.         Create a new cluster  $C_{new}$  based on  $X$ 
9.     endif
10. else  $\{N$  is a non-leaf node $\}$ 
11.      $C_{new} \leftarrow$  insert( $X, C_i.nextnode$ )
       $\{C_i.nextnode$  is the child node of CF entry  $C_i\}$ 
12.     Update statistics in  $C_i$ 
13. endif
14. if  $C_{new}$ 
       $\{a$  new CF entry needs to be inserted $\}$ 
15.     if node  $N$  has space for  $C_{new}$ 
16.         add  $C_{new}$  to  $N$ 
17.         return nil
18.     else
       $\{need$  to split  $N\}$ 
19.         create new node  $N_{new}$ , a and add  $C_{new}$ 
20.         seed  $N$  and  $N_{new}$  with the two most
      Distant clusters  $C_i$  and  $C_j$  in  $N$ 
21.         distribute  $C_{new}$  and the rest of  $C_k, k \neq i \neq j$ 
      According to their proximity to the seeds in  $N$  and  $N_{new}$ 
22.         return CF entry for  $N_{new}$ 
23.     endif
24. else
       $\{no$  new CF entry was created $\}$ 
25.     return nil
26. endif

```

รูปที่ 3.5 โค้ดเทียมสำหรับอัลกอริทึมเพิ่มเรคคอร์ด  $X$  ในซีเอฟทรี

### 3.3.2 การสรุปผล

เนื่องจากการแบ่งกลุ่มข้อมูลออกเป็นกลุ่มๆ แบบลำดับชั้นแบบซีเอฟทรีนั้นหมายความว่า โหนดที่เป็นโหนดพ่อจะเป็นโหนดที่รวบรวมข้อมูลสรุปของโหนดลูกๆ เอาไว้ อย่างไรก็ตามในการ

สรุปผลเราจำเป็นจะต้องเลือกกลุ่มข้อมูลที่น่าจะมีนัยสำคัญมากที่สุดในการแสดงผลและสรุปเป็นรายงานของระบบ

การเลือกกลุ่มข้อมูลที่มีนัยสำคัญนั้นขึ้นอยู่กับข้อมูลและจำนวนข้อมูลในแต่ละกลุ่มว่าสามารถบ่งบอกอะไรกับเราได้ อย่างเช่นการกระจายของข้อมูลภายในกลุ่มหรือขนาดของกลุ่มข้อมูลที่ใหญ่จะทำให้เราสามารถสังเกตเห็นได้ว่านั่นอาจจะเป็นนัยสำคัญอย่างหนึ่งเราจึงเลือกการสรุปผลของการแบ่งกลุ่มจากขนาดของกลุ่มข้อมูลที่ใหญ่จนผิดปกติซึ่งอาจบ่งบอกถึงการโจมตีในรูปแบบต่างๆได้ แต่อย่างไรก็ต้องมีเกณฑ์ซึ่งจะบอกได้ว่ากลุ่มข้อมูลใดใหญ่เกินไปเกณฑ์ที่ใช้คือค่า  $Tr$  ซึ่งเป็นเทรซโวล์ของขนาดกลุ่มข้อมูลหาได้จาก

$$Tr = r\tau \text{ ซึ่ง } r \in [0,1] \text{ และ } \tau \text{ แทนจำนวนเรคคอร์ดทั้งหมด}$$

ค่า  $r$  ซึ่งเป็นเทรซโวล์ของการแบ่งกลุ่มนั้นจะถูกกำหนดให้เป็น 0 ในช่วงเริ่มต้นซึ่งจะทำให้ทุกชุดข้อมูลที่เข้าผ่านกระบวนการวิเคราะห์ถูกจัดให้อยู่ในคนละกลุ่มซึ่งจะทำให้ต้องมีการแบ่งกลุ่มทุกๆครั้งที่มีเรคคอร์ดใหม่เข้ามา จนกระทั่งจำนวนหน่วยความจำที่ใช้จะเกินค่า  $M$  ซึ่งเป็นค่าหน่วยความจำที่มากที่สุดที่ระบบทำงานได้ ค่า  $r$  จะถูกทำให้เพิ่มขึ้นซึ่งมีผลทำให้จำนวนกลุ่มข้อมูลลดลงและใช้หน่วยความจำลดลงด้วย

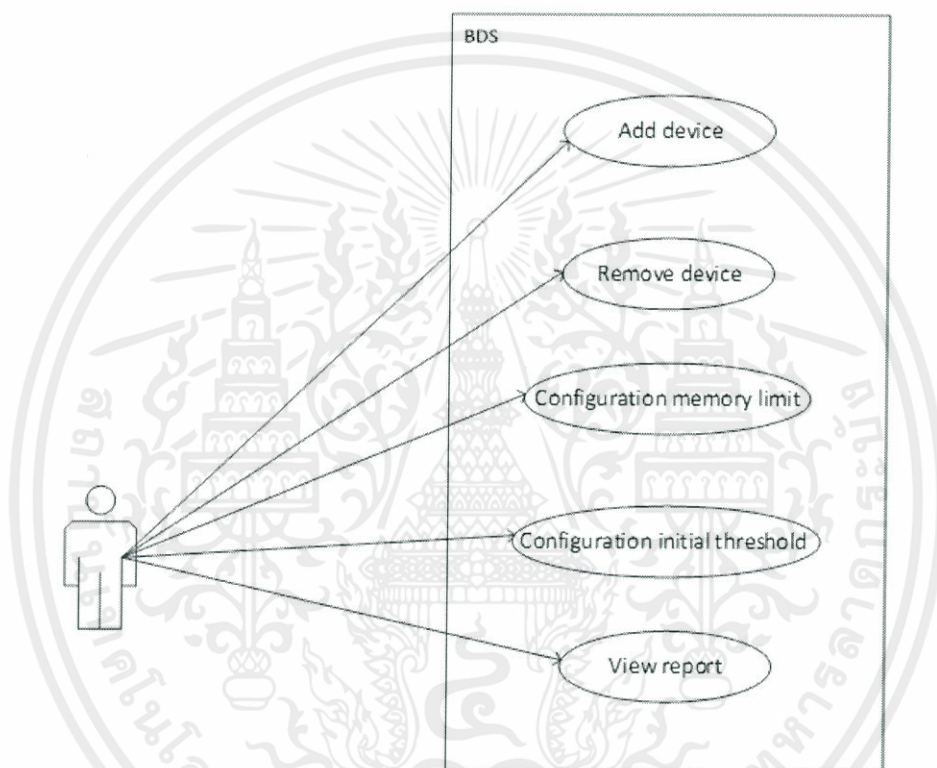
เมื่อได้กลุ่มข้อมูลที่มีขนาดข้อมูลมากพอที่จะนำมาวิเคราะห์นัยสำคัญ ระบบจะนำกลุ่มข้อมูลที่ได้นั้นมาทำการหาคาบการสื่อสารของแต่ละการเชื่อมต่อ ขึ้นอยู่กับว่าผู้ใช้ตั้งค่าระบบให้มีค่ากำหนดต่างๆเป็นเท่าใดอย่างเช่น มีการเชื่อมต่อทุกๆ 60 นาทีเป็นจำนวนมากกว่า 10 ครั้ง ข้อมูลเหล่านี้จะถูกนำไปพิจารณาเปรียบเทียบกับบอทที่มีข้อมูลอยู่อีกครั้งหนึ่ง หากการเปรียบเทียบนั้นมีระยะห่างของข้อมูลอยู่ในเกณฑ์ที่เหมาะสมข้อมูลของการเชื่อมต่อนั้นจะถูกกำหนดว่าเป็นข้อมูลที่ต้องสงสัย

การกำหนดการวิเคราะห์ข้อมูลขึ้นอยู่กับผู้ใช้งานต้องการให้ระบบทำการวิเคราะห์ข้อมูลเป็นช่วงเวลาเท่าใดโดยสามารถกำหนดได้เช่น เป็นช่วงวัน สัปดาห์ หรือแม้แต่ช่วงเดือน ซึ่งระบบจะสร้างรายงานเป็นกราฟต้นไม้และบทสรุปของรายงานนั้นเป็นไฟล์ TXT โดยที่สามารถเป็นอ่านได้จากระบบ

### 3.4 ส่วนของผู้ใช้งานและยูสเคสโตอะแกรม

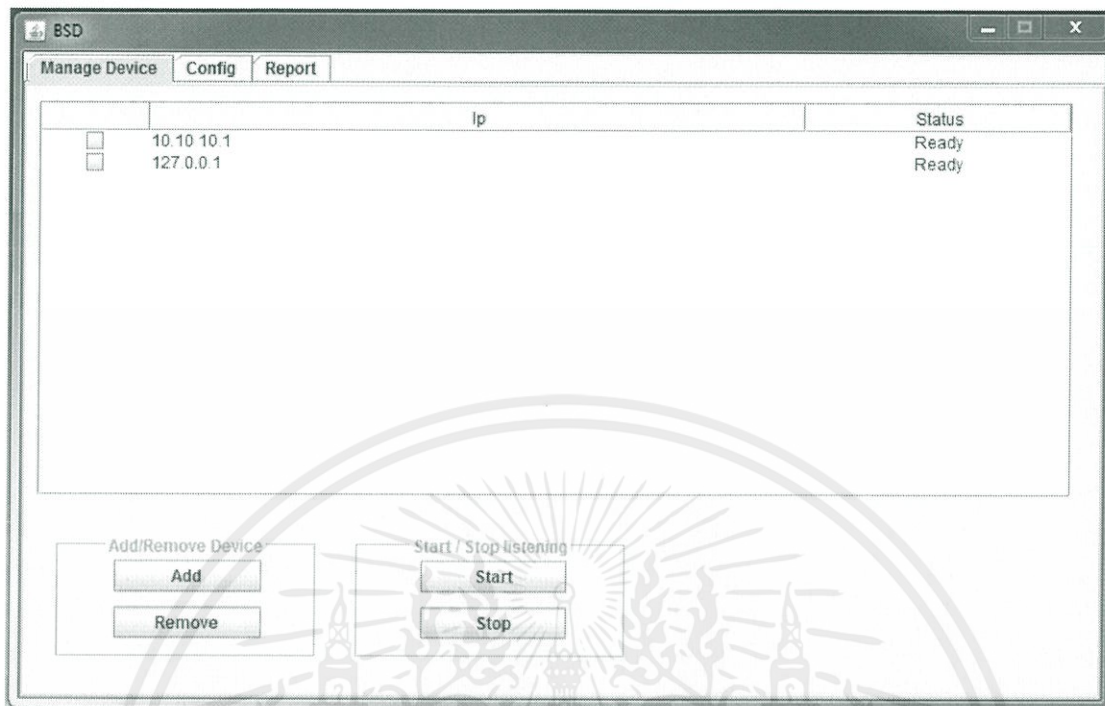
ส่วนของผู้ใช้นั้นสามารถตั้งค่าสำหรับการเพิ่มหมายเลขไอพีของเราเตอร์ สวิตช์ หรืออุปกรณ์ที่สามารถสร้างแพ็คเก็ตเน็ตโพล์ออกมาได้โดยสามารถเพิ่มและลบได้หลายอุปกรณ์โดยระบบจะรับข้อมูลจากอุปกรณ์เครือข่ายที่ได้ตั้งค่าไว้และเก็บลงฐานข้อมูลอัตโนมัติ ระบบจะทำการสร้างและจับคู่

อุปกรณ์กับฐานข้อมูลให้โดยอัตโนมัติจึงทำให้ผู้ใช้ไม่จำเป็นต้องกังวลเกี่ยวกับการจัดการเกี่ยวกับฐานข้อมูล นอกจากนี้ผู้ใ้ยังสามารถเริ่มและหยุดการรับข้อมูลจากอุปกรณ์เครือข่ายได้ตามต้องการ ผู้ใช้จำเป็นต้องมีการกำหนดค่าเริ่มต้นบางอย่างให้กับการรับข้อมูลและจัดการเกี่ยวกับซีเอฟทีรีในเบื้องต้นก่อนคือการตั้งค่าหน่วยความจำที่สามารถใช้ได้หน่วยเป็นเมกะไบต์และต้องตั้งค่าเทรชโอล์สำหรับบรัคมีของแต่ละกลุ่มข้อมูลก่อนโดยค่าปกติจะถูกตั้งเป็น 0

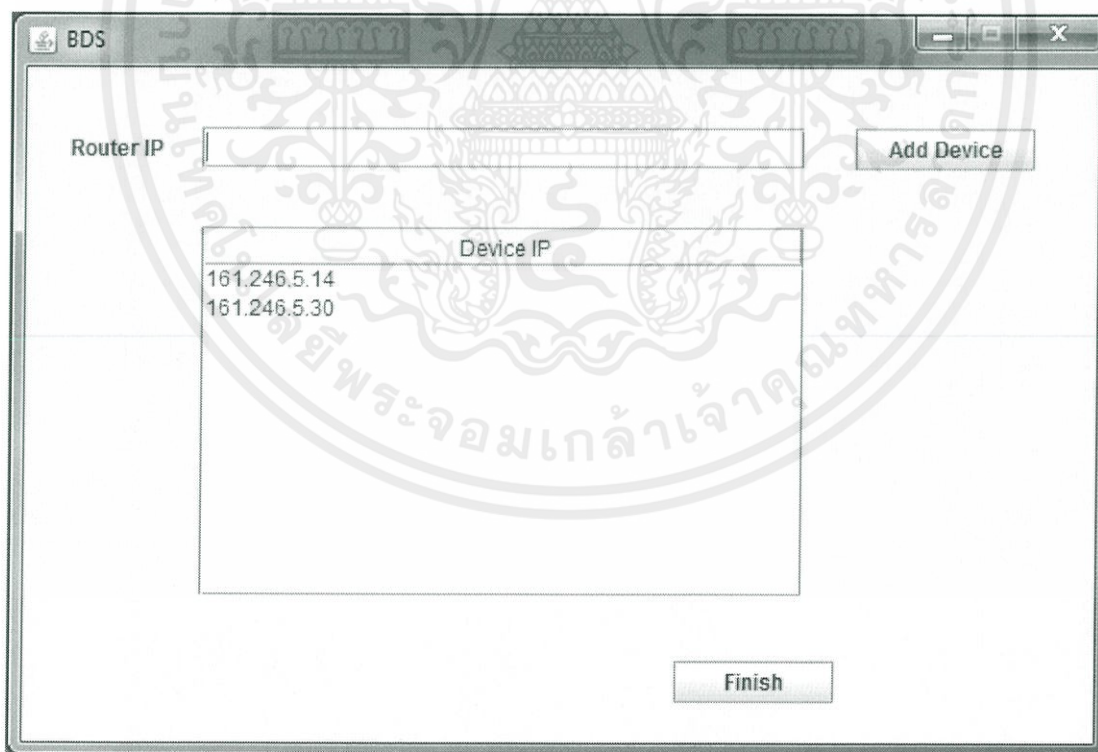


รูปที่ 3.6 ยูสเคสไดอะแกรม

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

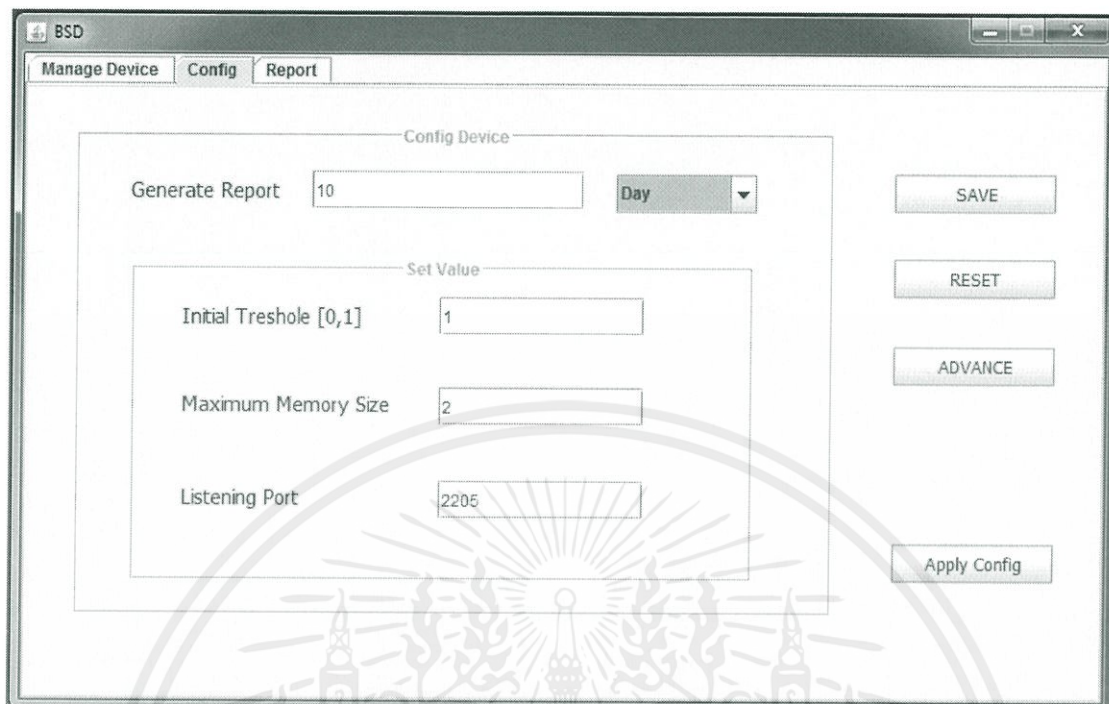


รูปที่ 3.7 หน้าหลักของระบบซึ่งเป็นหน้าสำหรับการจัดการอุปกรณ์

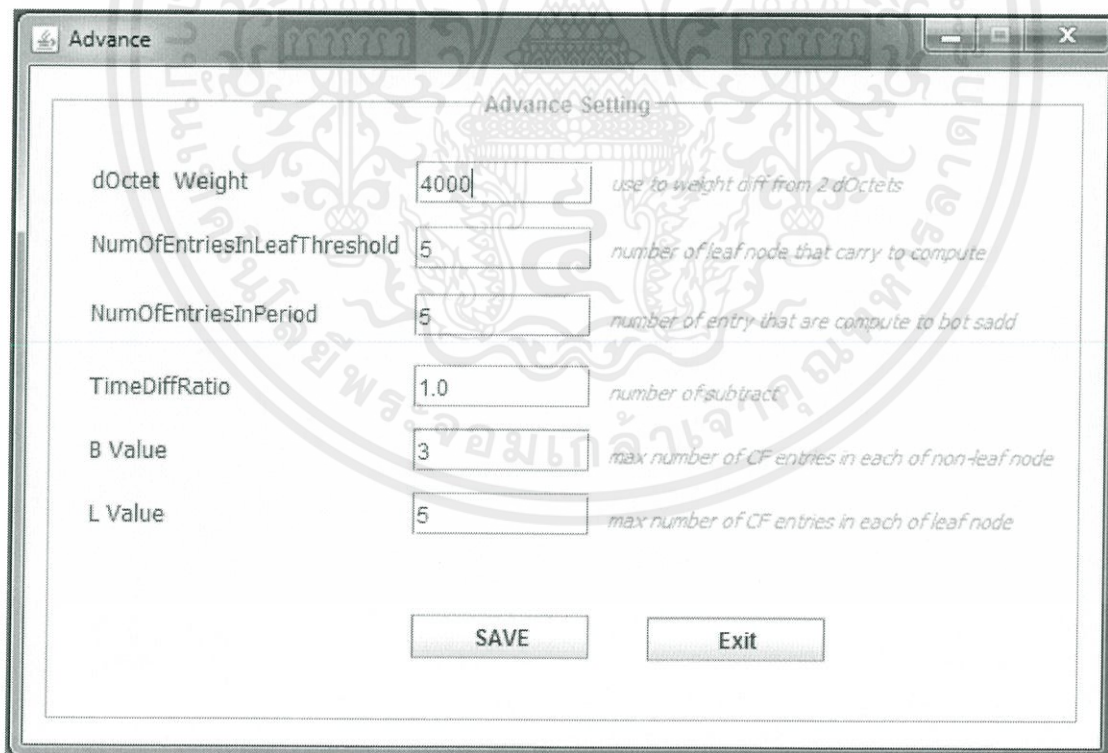


รูปที่ 3.8 หน้าเพิ่มไอพีแอดเดรสของเราเตอร์ให้กับโปรแกรม

เอกสารนี้เป็นเอกสารที่สงวนลิขสิทธิ์ของมหาวิทยาลัยเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 3.9 หน้าสำหรับตั้งค่าอุปกรณ์



รูปที่ 3.10 หน้าสำหรับตั้งค่าอุปกรณ์ขั้นสูง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับใช้ภายในเท่านั้น ไม่ควรเผยแพร่ให้บุคคลภายนอก  
 ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

BSD

Manage Device Config Report

View Report Now!

Select Report: 2014-02-17T05-30-16.txt

View Custom

Custom will take a long time because it will re...

Search Connection

Source IP:  Port:

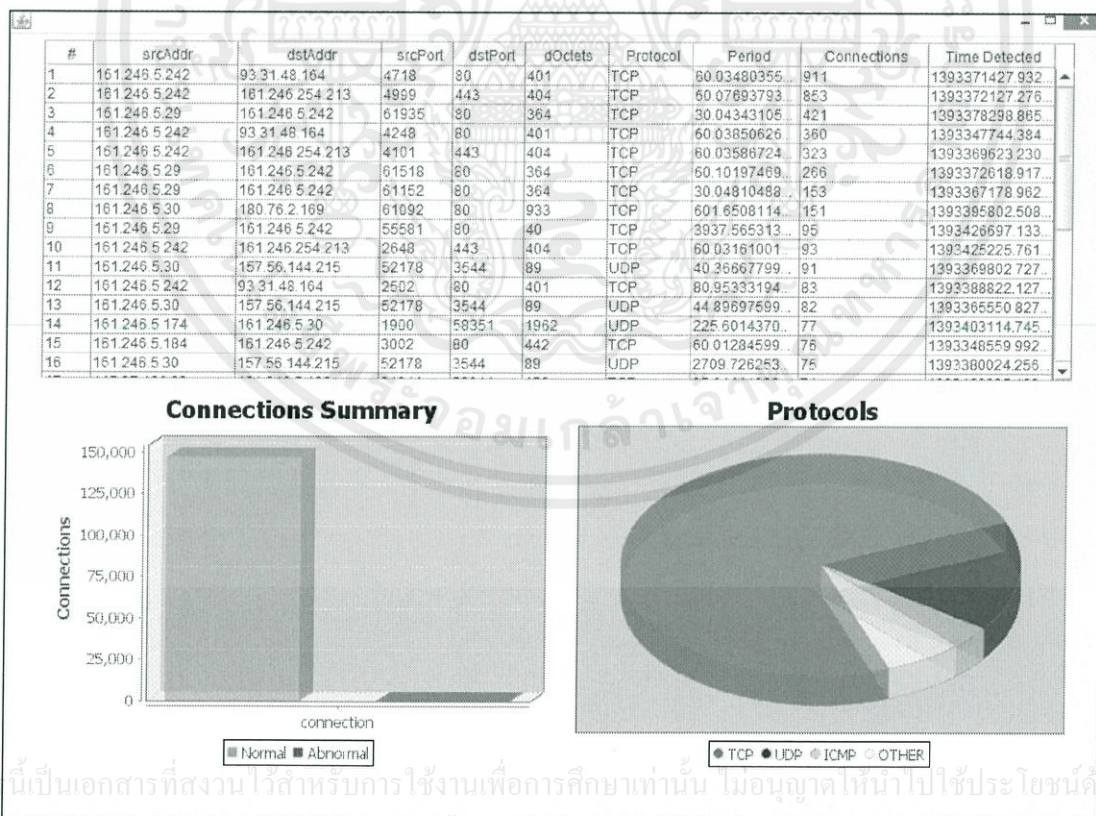
Destination IP:  Port:

dOctets:  Byte:

Period:  Secs:

Search

รูปที่ 3.11 หน้าการตั้งค่าการแสดงผล



รูปที่ 3.12 หน้ารายงานจำนวนการเชื่อมต่อและกราฟแสดงผลจำนวนบอทเน็ตที่มีในเครือข่าย

## บทที่ 4

### การทดลองและผลการทดลอง

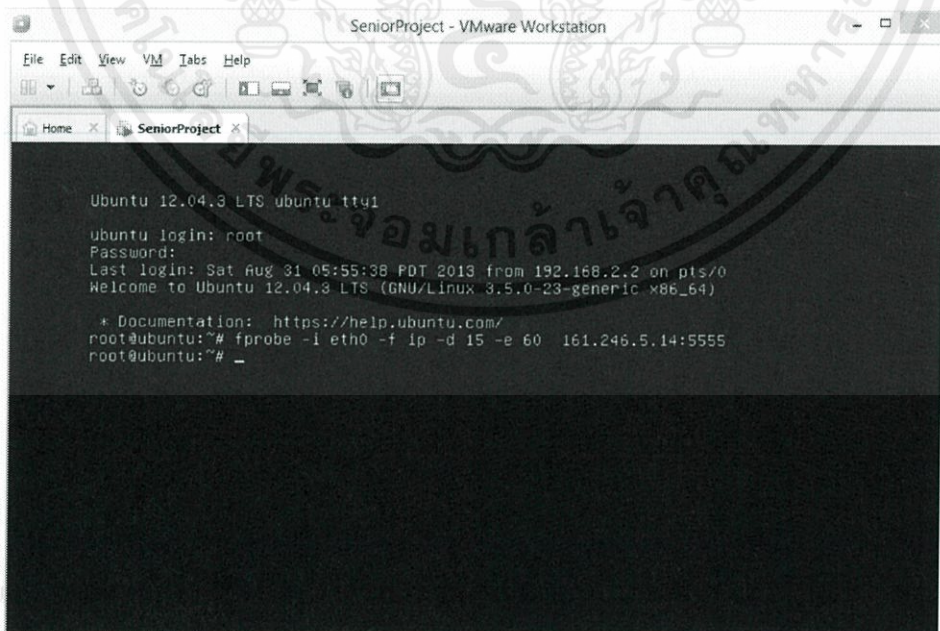
#### 4.1 การทดลองโปรแกรมเอฟโพรบ (Fprobe)

ในการมอนิเตอร์เครือข่ายโดยใช้แพ็คเก็ตเน็ตโพล์ ที่ซิสโก้ ได้พัฒนาขึ้นมาตั้งแต่ใช้ อุปกรณ์เครือข่ายที่รองรับการใช้งานของ เน็ตโพล์ แต่เนื่องจากสาขาวิชาวิศวกรรมคอมพิวเตอร์ยังไม่ มีอุปกรณ์ที่สามารถรองรับการทำงานของโปรโตคอล เน็ตโพล์ ได้ ผู้ดำเนินโครงการจึงจำเป็นต้องหา วิธีที่จะมอนิเตอร์เครือข่ายได้อย่างมีประสิทธิภาพ

เอฟโพรบ เป็นซอฟต์แวร์หนึ่งที่มีความสามารถในการสร้างแพ็คเก็ตเน็ตโพล์ จากการแพ็คเก็ต ต่างๆที่วิ่งผ่านเครื่องที่ลงซอฟต์แวร์นี้ โดยที่ เอฟโพรบ รองรับการทำงานของ เน็ตโพล์ เวอร์ชันที่ 5 ซึ่งมีรายละเอียดและข้อมูลต่างๆที่ผู้ดำเนินโครงการต้องการ

##### 4.1.1 การติดตั้งและตั้งค่าเอฟโพรบ (Fprobe) และการเก็บข้อมูล NetFlow

การทำงานของ เอฟโพรบ นั้นต้องทำงานบนระบบปฏิบัติการ Linux เท่านั้นซึ่งผู้ทดลองได้ทำการติดตั้งระบบปฏิบัติการ Ubuntu Server เวอร์ชัน 12.04.3 LTS บนเวอร์ชวลแมชีน VMware เวอร์ชัน 9

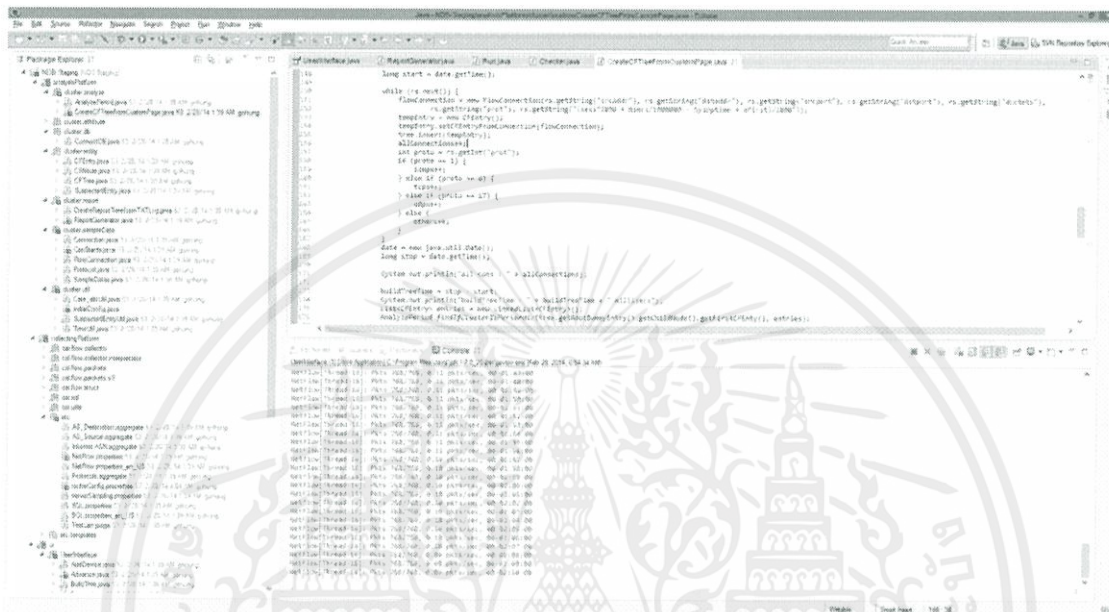


```
SeniorProject - VMware Workstation
File Edit View VM Tabs Help
Home x SeniorProject x
Ubuntu 12.04.3 LTS ubuntu tty1
ubuntu login: root
Password:
Last login: Sat Aug 31 05:55:38 PDT 2013 from 192.168.2.2 on pts/0
Welcome to Ubuntu 12.04.3 LTS (GNU/Linux 3.5.0-23-generic x86_64)

* Documentation: https://help.ubuntu.com/
root@ubuntu:~# fprobe -i eth0 -f ip -d 15 -e 60 161.246.5.14:5555
root@ubuntu:~#
```

รูปที่ 4.1 แสดงการรันคำสั่ง `fprobe -i eth0 ip -d 15 -e 60 161.246.5.14:5555`

เพื่อสั่งให้ซอฟต์แวร์ทำการเก็บข้อมูลและสร้างเป็นแพ็คเกจเน็ตเวิร์กส่งไปที่เครื่องคอมพิวเตอร์ของผู้วิจัย ผู้วิจัยได้ทำการพัฒนาโปรแกรมส่วนของการรวบรวมข้อมูล NetFlow ที่ได้รับมาจากเครื่องเซิร์ฟเวอร์ เพื่อกรองเฉพาะข้อมูลที่จำเป็นเพื่อเก็บข้อมูลลงฐานข้อมูล



รูปที่ 4.2 รูปภาพแสดงการทดลองเก็บข้อมูล NetFlow

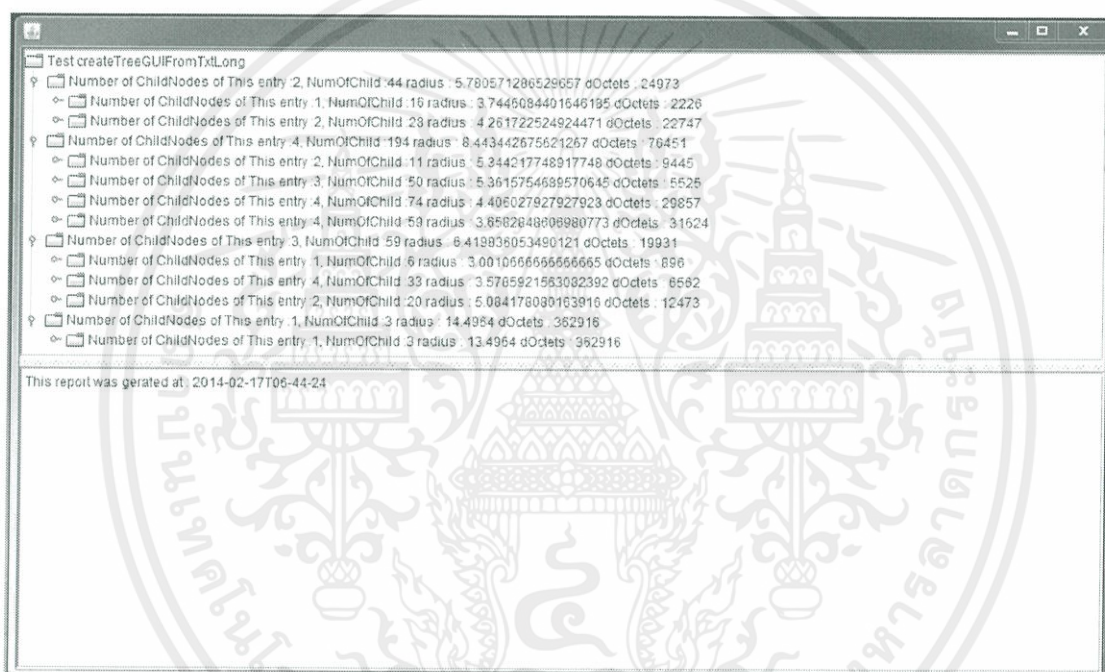
Router IP	SysUptime	Srcs	Dests	Flow_Sequene	Engine_Type	Engine_ID	SrcAddr	DestAddr	NextHop	Inpuot	Output	dPkts	dOctets	dPst	sLast	Src-Port	De-Port
10.0.0.1	1048509955	1333340059	1562000	663106	0	0	161.246.5.184	161.246.5.242	0.0.0.0	0	0	6	442	1048550127	1048558147	1151	80
10.0.0.1	1048509955	1333340059	1562000	663106	0	0	161.246.5.242	161.246.5.184	0.0.0.0	0	0	4	1762	1048558151	1048558164	80	1152
10.0.0.1	1048509955	1333340059	1562000	663106	0	0	161.246.5.184	161.246.5.242	0.0.0.0	0	0	5	402	1048558150	1048558155	1152	80
10.0.0.1	1048509955	1333340059	1562000	663106	0	0	161.246.5.225	161.246.5.30	0.0.0.0	0	0	13	19555	1048558040	1048558055	8157	52750
10.0.0.1	1048509955	1333340059	1562000	663106	0	0	161.246.5.225	224.0.0.22	0.0.0.0	0	0	6	248	1048558028	1048559091	0	0
10.0.0.1	1048570000	1333340059	4903000	663136	0	0	161.246.5.30	161.246.5.225	0.0.0.0	0	0	13	1491	1048556024	1048556056	52751	5357
10.0.0.1	1048570000	1333340059	4903000	663136	0	0	161.246.5.184	161.246.5.242	0.0.0.0	0	0	6	442	1048558161	1048558171	1153	80
10.0.0.1	1048570000	1333340059	4903000	663136	0	0	161.246.5.242	161.246.5.242	0.0.0.0	0	0	4	301	1048558257	1048558321	443	1413
10.0.0.1	1048570000	1333340059	4903000	663136	0	0	161.246.5.225	161.246.5.30	0.0.0.0	0	0	2	2544	1048556043	1048556156	3702	61122
10.0.0.1	1048570000	1333340059	4903000	663136	0	0	161.246.5.242	161.246.5.184	0.0.0.0	0	0	4	1762	1048558175	1048558185	80	1154
10.0.0.1	1048570000	1333340059	4903000	663136	0	0	161.246.5.242	161.246.5.242	0.0.0.0	0	0	4	301	1048558266	1048582817	443	1411
10.0.0.1	1048570000	1333340059	4903000	663136	0	0	161.246.5.242	161.246.254.213	0.0.0.0	0	0	6	454	1048558269	1048582826	1411	443
10.0.0.1	1048570000	1333340059	4903000	663136	0	0	161.246.5.242	93.31.48.164	0.0.0.0	0	0	5	401	1048558269	1048582844	1410	80
10.0.0.1	1048580997	1333340064	1867000	663156	0	0	161.246.5.184	161.246.5.242	0.0.0.0	0	0	6	442	1048558174	1048558186	1154	80
10.0.0.1	1048580997	1333340064	1867000	663156	0	0	93.31.48.164	161.246.5.242	0.0.0.0	0	0	5	288	1048558259	1048558265	80	1410
10.0.0.1	1048580997	1333340064	1867000	663156	0	0	161.246.5.242	93.31.48.164	0.0.0.0	0	0	5	401	1048558287	1048558294	1412	80
10.0.0.1	1048580997	1333340064	1867000	663156	0	0	161.246.5.225	161.246.5.30	0.0.0.0	0	0	12	19555	1048558046	1048558066	5157	52751
10.0.0.1	1048580997	1333340064	1867000	663156	0	0	161.246.5.242	94.245.121.251	0.0.0.0	0	0	6	454	1048558286	1048558321	1413	443
10.0.0.1	1048580997	1333340064	1867000	663156	0	0	161.246.5.30	94.245.121.251	0.0.0.0	0	0	1	89	1048558491	1048558491	6316	2644
10.0.0.1	1048580997	1333340064	1867000	663156	0	0	161.246.5.242	161.246.5.184	0.0.0.0	0	0	4	1762	1048558161	1048558170	80	1153
10.0.0.1	1048580997	1333340064	1867000	663156	0	0	93.31.48.164	161.246.5.242	0.0.0.0	0	0	5	268	1048558269	1048558296	80	1412
10.0.0.1	1048580997	1333340064	1867000	663156	0	0	161.246.5.225	161.246.5.30	0.0.0.0	0	0	2	2544	1048558040	1048558014	3702	49564
10.0.0.1	1048570000	1333340064	7545000	663194	0	0	94.245.121.251	161.246.5.30	0.0.0.0	0	0	1	137	1048558770	1048558770	3844	63316
10.0.0.1	1048570000	1333340064	7545000	663194	0	0	161.246.5.242	93.31.48.164	0.0.0.0	0	0	5	401	1048558354	1048558359	1416	80
10.0.0.1	1048570000	1333340064	7545000	663194	0	0	161.246.5.242	161.246.254.213	0.0.0.0	0	0	6	454	1048558326	1048558354	1415	443
10.0.0.1	1048570000	1333340064	7545000	663194	0	0	93.31.48.164	161.246.5.242	0.0.0.0	0	0	5	268	1048558321	1048558325	80	1414

รูปที่ 4.3 รูปภาพแสดงฐานข้อมูล NetFlow

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดลอกเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## 4.2 การทดลองการทำงานของอัลกอริทึมเบิร์ช

การทดลองการทำงานของอัลกอริทึมเบิร์ช ผู้วิจัยได้พัฒนาส่วนของการแบ่งกลุ่มข้อมูลโดยยึดหลักการตามเอกสารอ้างอิง ในขั้นแรกของการแบ่งกลุ่มข้อมูลจะต้องมีการเชื่อมต่อกับฐานข้อมูลเพื่อร้องขอข้อมูลดิบที่จะนำมาวิเคราะห์ จากนั้นจึงทำการแปลงข้อมูลดิบเป็นรูปแบบสำหรับการวิเคราะห์ แล้วจึงแบ่งกลุ่มข้อมูล ซึ่งหลังจากการทดลองแบ่งกลุ่มข้อมูลเครือข่ายด้วยอัลกอริทึมเบิร์ชแล้วพบว่า เวลาที่ใช้ในการร้องขอข้อมูลจากฐานข้อมูลนั้นมักจะมีอัตราส่วนมากกว่าส่วนคำนวณเพื่อแบ่งกลุ่มข้อมูล



รูปที่ 4.4 ผลลัพธ์การแบ่งกลุ่มตามอัลกอริทึม

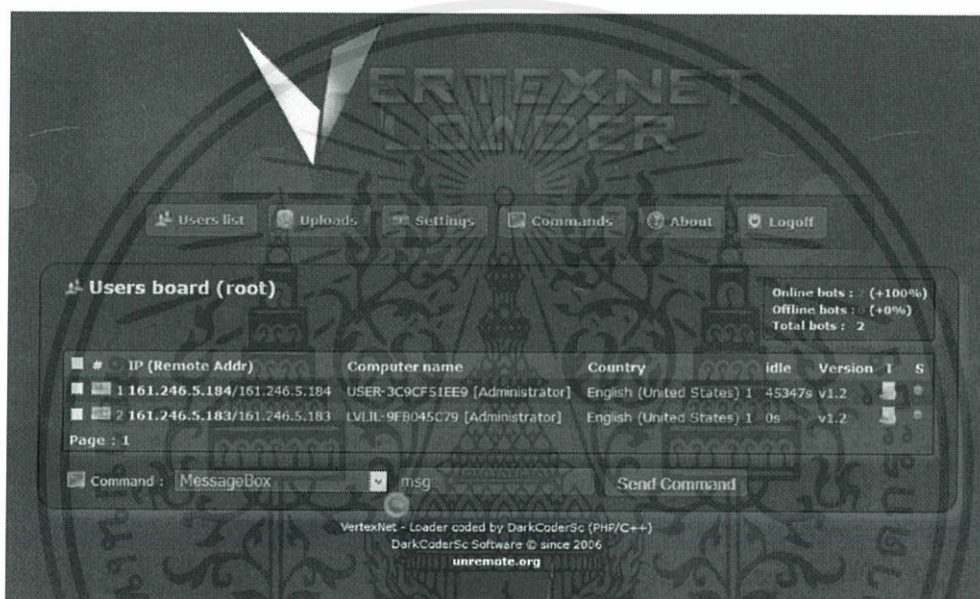
## 4.3 การทดลองการทำงานของ Botnet (VertexNet รุ่น 1.2.1)

### 4.3.1 ทดลองการติดตั้งบอท

ในการทดลองการทำงานของ Botnet นั้นได้เลือกใช้ Botnet ชื่อ VertexNet รุ่น 1.2.1 ซึ่งเป็น Botnet ที่สามารถหาใช้งานได้ง่ายจากการดาวน์โหลดทางอินเทอร์เน็ต VertexNet นั้นประกอบไปด้วยสองส่วนคือส่วนที่เป็น Command and Control Server (C&C) ซึ่งเป็นส่วนที่ทำงานเป็นตัวสั่งงานบอทและ Loader ซึ่งทำหน้าที่เป็นส่วนแพร่กระจายบอท

การทดลองการทำงานของบอทนั้นได้ถูกจัดการให้มามีการทำงานอยู่ในสิ่งแวดล้อมเสมือน โดยติดตั้ง C&C บนเวอร์ชวลแมชีนโดยใช้ระบบปฏิบัติการ XP มีหน่วยความจำหลัก 512 เมกะไบต์และ

หน่วยความจำรอง 30 กิกะไบต์และได้จำกัดการเชื่อมต่อสู่เครือข่ายภายนอกของเวอร์ชวลแมชีนเพื่อเป็นการควบคุมไม่ให้ส่งผลกระทบต่อเครือข่ายภายนอก C&C จะทำงานเป็น HTTP เซิร์ฟเวอร์สำหรับสั่งงานและรับการเชื่อมต่อจากบอท โดย C&C ถูกติดตั้งโดยกำหนดหมายเลขไอพี 161.246.5.182/24 และได้ทำการติดตั้งบอทจำนวน 2 เครื่องคือหมายเลขไอพี 161.246.5.183/24 และ 161.246.5.184/24 ซึ่งผลของการทดลองการเชื่อมต่อระหว่างบอทและ C&C คือตัว C&C เซิร์ฟเวอร์สามารถติดต่อกับบอทได้นอกจากนั้นยังสามารถรับรู้ได้ว่าบอทนั้นยังอยู่บนเครือข่ายหรือไม่ดังรูป



รูปที่ 4.5 การเชื่อมต่อระหว่าง C&C กับบอทโดยที่จะมีค่า “idle” ซึ่งเป็นค่าที่บอกว่าบอทนั้นมีการเคลื่อนไหวหรือไม่ ถ้าหากมีการเคลื่อนไหวของบอทเช่นผู้ใช้งานมีการขยับเมาส์บอทก็จะส่งแพ็คเกจเพื่อรีเซ็ตค่า “idle” เป็น 0

#### 4.3.2 ทดลองการสั่งงานบอทและการสังเกตผล

ในการสั่งงานบอทของ C&C เซิร์ฟเวอร์นั้นจะมีคำสั่งที่ได้ถูกกำหนดไว้ให้แล้วว่าจะต้องใช้คำสั่งแบบใดโดยคำสั่งที่ให้ใช้นั้นมีอยู่ทั้งหมด 14 คำสั่งซึ่งประกอบไปด้วย

msg : การส่งข้อความ

exec : การรันคำสั่ง (Command Prompt) ของเครื่องที่เป็นบอท

close : ปิดการทำงานของ loader ของเครื่องที่เป็นบอท

urldl : ดาวน์โหลดไฟล์จากเครื่องที่เป็นบอท

getproc : คำสั่งเพื่อดูโปรเซสที่เปิดทำงานอยู่ในเครื่องที่เป็นบอท

setkeylogger : เปิดการทำงานของ key logger

getlogs : อ่านไฟล์บันทึกของ key logger

readfile : อ่านไฟล์ที่เครื่องที่เป็นบอท

uninstall : ถอนการติดตั้ง loader ออกจากเครื่องที่เป็นบอท

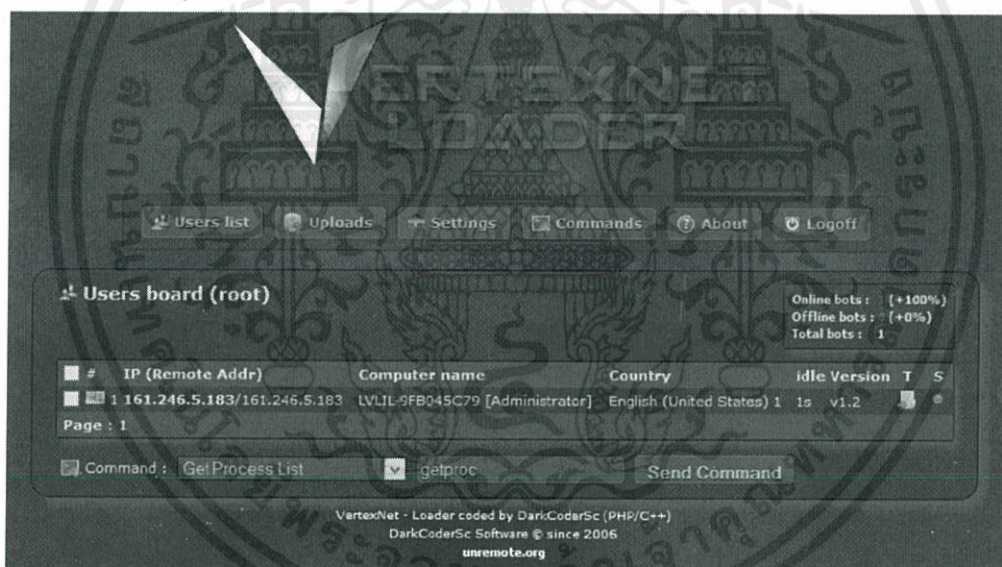
httpflood : สั่งให้บอททำการโจมตีแบบ httpflood ไปยังเครื่องเป้าหมาย

remoteshell : ทำการรีโมทใช้งานเซสชันที่สามารถใช้งานบนเครื่องบอทที่ใช้งานลินุกซ์

visitpage : สั่งงานให้เครื่องที่เป็นบอททำการเปิดเว็บเพจที่ต้องการช่วงขณะหนึ่ง

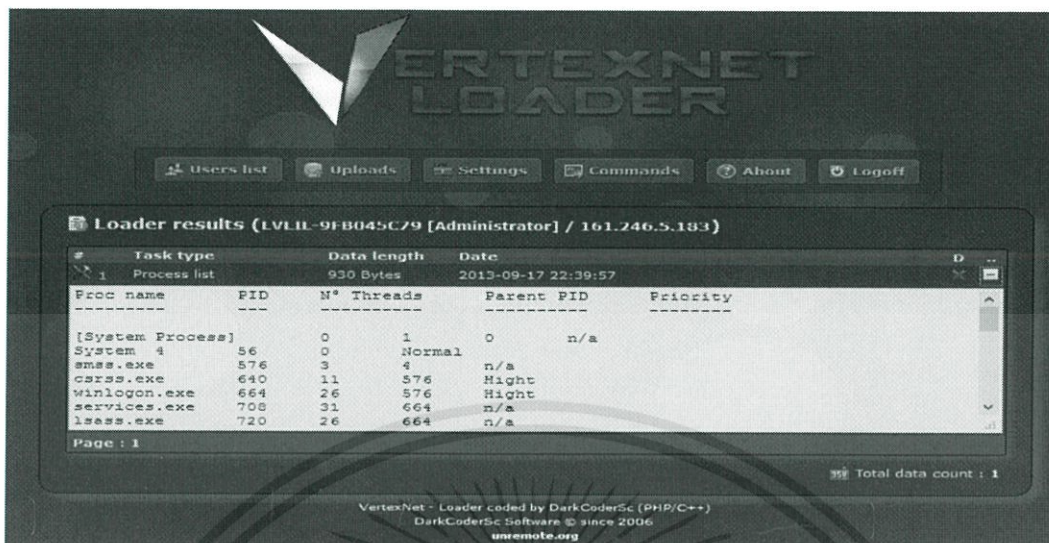
update : คำสั่งอัปเดต loader ให้เป็นเวอร์ชันล่าสุด

การทดลองการสั่งงานของ C&C ได้ทดลองใช้งานคำสั่ง “getproc” และคำสั่ง “exec : shutdown -r -t 0” ซึ่งคำสั่ง getproc นั้นจะมีผลทำให้ได้รับการรีเทิร์นโพเซสที่ทำงานอยู่บนเครื่องที่เป็นบอท และคำสั่ง shutdown เป็นคำสั่งซึ่งสั่งให้เครื่องบอททำการเริ่มการทำงานใหม่ ซึ่งทั้งสองคำสั่งมีผลในทันทีที่ C&C สามารถติดต่อกับบอทได้ใช้เวลาประมาณ 3 -5 วินาที

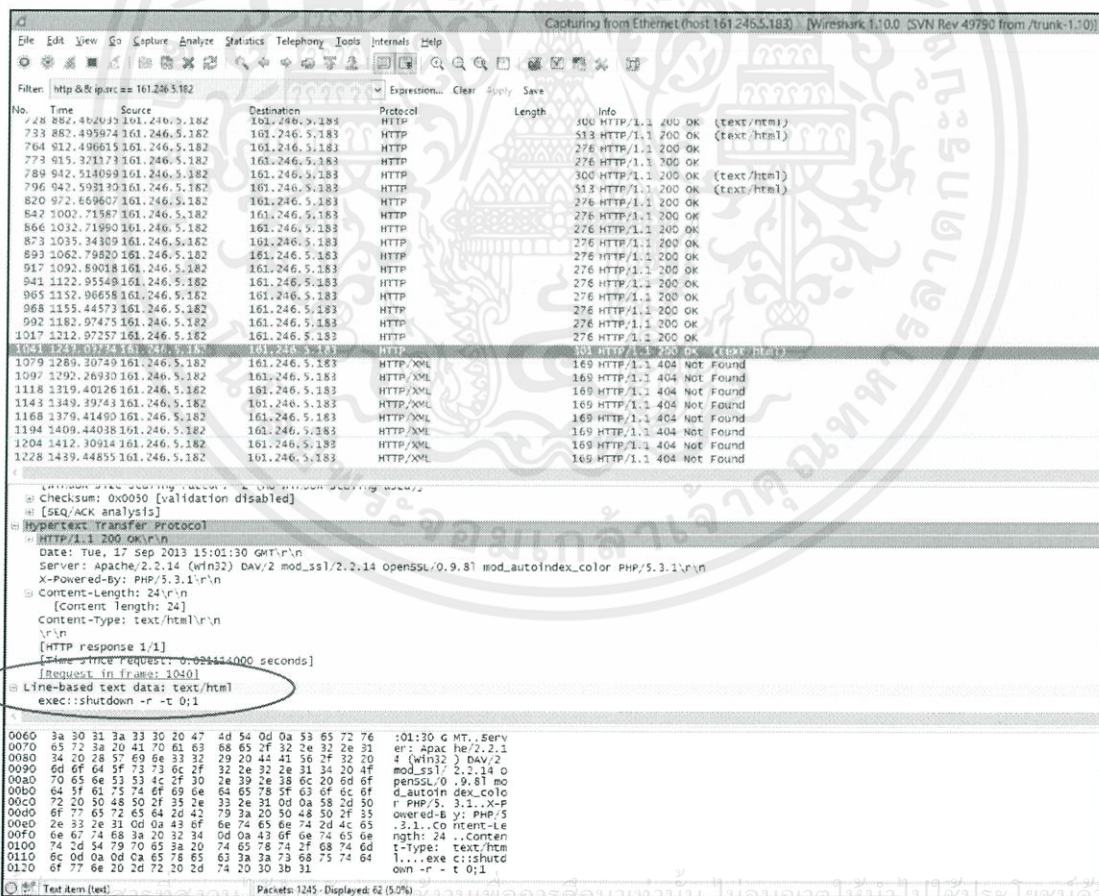


รูปที่ 4.6 การทดลองสั่งงานบอทด้วยคำสั่ง getproc

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.7 ผลลัพธ์ของคำสั่ง getproc ซึ่งจะเห็นได้ว่ามีการ  
ริเทิร์นโปรเซสที่ทำงานอยู่บนเครื่องที่เป็นบอท



เอกสาร

ไม่ว่ากรณีใดๆ ทั้งสิ้น รูปที่ 4.8 การทดลองการดักจับแพ็คเก็ตเพื่อสังเกตการส่งงานของ C&C โดย  
เป็นการดักจับแพ็คเก็ตเกิดขณะที่มีการส่งคำสั่ง shutdown -r -t 0

#### 4.4 การทดลองการทำงานของ Botnet (Zeus รุ่น 2.1.0.1)

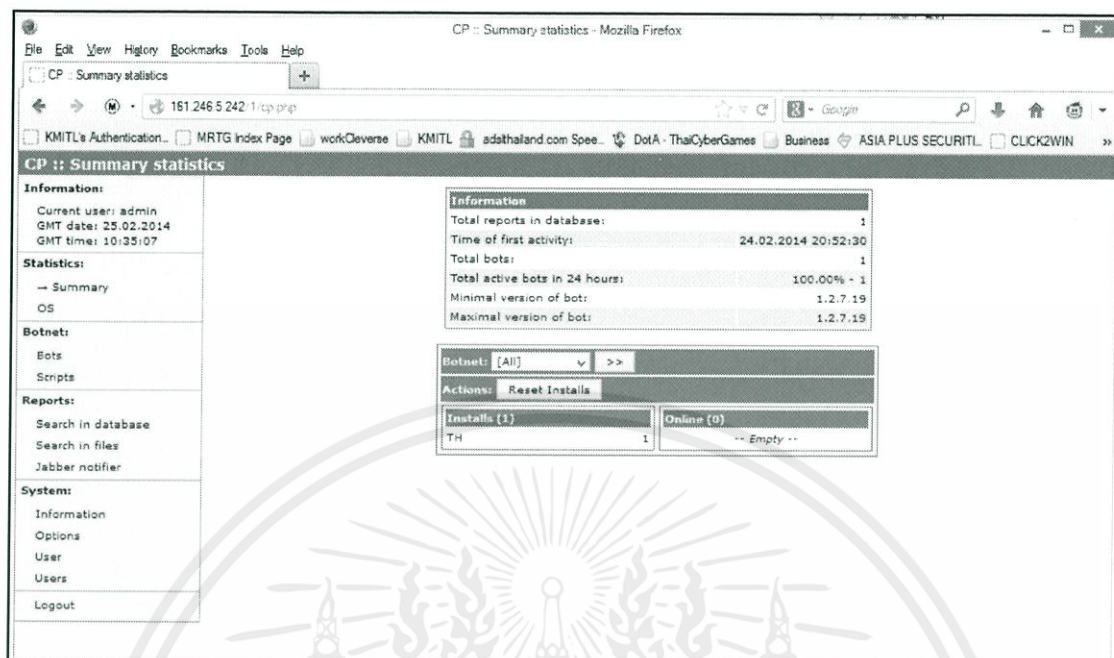
การทดลองการใช้งานบอท Zeus รุ่น 2.1.0.1 นั้น ผู้วิจัยสามารถดักจับแพ็คเกจของ Zeus ได้ และจากผลการทดลองวิเคราะห์เห็นว่า การทำงานของ Zeus มีลักษณะเป็นคานการทำงานซ้ำๆตาม ทฤษฎีที่ได้อ้างถึงก่อนหน้า ทำให้ผู้วิจัยสามารถพัฒนาระบบเพื่อวิเคราะห์ตรวจหาคานการทำงานได้ นอกจากนี้ยังเห็นว่าการทำงานของ Zeus นั้นจะมีการร้องขอไปยัง HTTP Server โดยมีพารามิเตอร์ ข้างเดิมคือร้องขอหน้า page /cfg.bin เสมอซึ่งแสดงให้เห็นว่ามีความเป็นไปได้ที่จะสามารถพัฒนา IDS เพื่อตรวจจับลักษณะดังกล่าวได้

The screenshot shows a Wireshark capture of network traffic. The top pane displays a list of captured packets, all of which are HTTP GET requests for the resource /cfg.bin. The bottom pane shows the detailed view of a selected packet, including the Ethernet II header, Internet Protocol Version 4 header, and the Hypertext Transfer Protocol (HTTP) section, which shows a GET request for /cfg.bin with various headers like User-Agent, Host, and Cache-Control.

No.	Time	Source	Destination	Protocol	Source Port	Destination Port	Host	Total Length	Info
59287	36622.24000	161.246.5.184	161.246.5.242	HTTP	1154	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
60690	7232.27068	161.246.5.184	161.246.5.242	HTTP	2252	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
22961	10832.2816	161.246.5.184	161.246.5.242	HTTP	4115	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
83254	14432.3181	161.246.5.184	161.246.5.242	HTTP	1074	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
50030	18032.3460	161.246.5.184	161.246.5.242	HTTP	1987	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
19618	21632.3608	161.246.5.184	161.246.5.242	HTTP	2891	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
02933	25232.3794	161.246.5.184	161.246.5.242	HTTP	3816	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
03253	28832.4108	161.246.5.184	161.246.5.242	HTTP	4795	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
18059	32432.4246	161.246.5.184	161.246.5.242	HTTP	1724	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
26140	36032.4730	161.246.5.184	161.246.5.242	HTTP	2659	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
46872	39632.4941	161.246.5.184	161.246.5.242	HTTP	3578	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1

รูปที่ 4.9 ผลที่ได้จากการดักจับแพ็คเกจของ Wireshark ที่เป็นของบอท Zeus ที่ร้องขอไปที่ /cfg.bin

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.10 หน้าเว็บของบอทมาสเตอร์ที่คอยควบคุมและส่งคำสั่งให้กับบอทลูก

อย่างไรก็ตามการเชื่อมต่อของบอทกับเซิร์ฟเวอร์นั้นจะมีลักษณะซ้ำกันในคาบเวลาที่ซ้ำกันเสมอ แต่ผู้ใช้งานบอทก็สามารถที่จะตั้งค่าการทำงานของบอทให้มีคาบเวลาที่ต่างกันไปได้ซึ่งทำให้ยากต่อการตรวจจับ ดังนั้นการตรวจจับบอทโดยใช้พฤติกรรมทางเครือข่ายสามารถตรวจจับได้จากเพียงการค้นหากการทำงานเป็นคาบๆ ของการเชื่อมต่อใดๆ เท่านั้น

#### 4.4 การทดลองการวิเคราะห์การสื่อสารของบอทกับเซิร์ฟเวอร์

##### 4.4.1 พฤติกรรมการเชื่อมต่อเป็นคาบของบอทเน็ต Zeus รุ่น 2.1.0.1 และ VertexNet รุ่น 1.2.1

ลักษณะการสร้างการเชื่อมต่อเป็นช่วงๆ โดยที่มีการส่งข้อมูลที่มีขนาดเท่าๆกันไปยังหมายเลขไอพีปลายทางซ้ำกันหลายๆครั้งทำให้เป็นที่น่าสังเกตว่าพฤติกรรมดังกล่าวอาจไม่ใช่การส่งข้อมูลปกติ โดยจากการทดลองการทำงานของบอทเน็ต Zeus โดยตั้งค่าให้มีการส่งข้อมูลเพื่อติดต่อกับเซิร์ฟเวอร์เป็นคาบมาตรฐานของบอทรุ่นนี้แล้วอาศัยเครื่องมือดักจับแพ็คเกจ Wireshark ดักจับแพ็คเกจเป็นเวลา 12 ชั่วโมงนับแต่เครื่องเหยื่อติดบอท) ของการส่งข้อมูลพบว่าบอท Zeus มีการส่งข้อมูลเป็นคาบอย่างมีนัยสำคัญคือการส่งข้อมูลติดต่อกับเครื่องเซิร์ฟเวอร์เป็นคาบเวลา 60 นาทีโดยมีขนาดข้อมูลที่เท่ากันเสมอคือ 207 ไบต์ ติดต่อกัน โดยที่เมื่อเปรียบเทียบกับกรส่งข้อมูลปกติของผู้ใช้งานจะเห็นความแตกต่างอย่างเห็นได้ชัด

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการศึกษาเท่านั้น ไม่นิยามให้นำไปใช้ประโยชน์ในการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งยังมีให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

No.	Time	Source	Destination	Protocol	Source Port	Destination Port	Host	Total Length	Info
47	23202.217	161.246.5.184	161.246.5.242	HTTP	1358	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
96287	3632.240	161.246.5.184	161.246.5.242	HTTP	2253	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
60690	7232.2706	161.246.5.184	161.246.5.242	HTTP	3184	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
22961	10832.2816	161.246.5.184	161.246.5.242	HTTP	4115	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
83254	14432.3181	161.246.5.184	161.246.5.242	HTTP	1074	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
50030	18032.3460	161.246.5.184	161.246.5.242	HTTP	1987	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
19618	21632.3608	161.246.5.184	161.246.5.242	HTTP	2891	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
02933	25232.3794	161.246.5.184	161.246.5.242	HTTP	3816	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
03253	28832.4108	161.246.5.184	161.246.5.242	HTTP	4758	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
18059	32432.4246	161.246.5.184	161.246.5.242	HTTP	1724	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
26140	36032.4730	161.246.5.184	161.246.5.242	HTTP	2659	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1
46872	39632.4941	161.246.5.184	161.246.5.242	HTTP	3578	80	161.246.5.242	203	GET /1/cfg.bin HTTP/1.1

รูปที่ 4.11 รูปแสดงการกรองแพ็คเก็ตจากหมายเลขไอพี, หมายเลขพอร์ต, โพรโทคอลและปริมาณข้อมูล

จากการทดลองดักจับแพ็คเก็ตของเครื่องที่ถูกคุกคามโดย VertexNet รุ่น 1.2.1 พบว่าเครื่องที่ถูกคุกคามนั้นมีการสร้างการเชื่อมต่อไปยังเครื่องเซิร์ฟเวอร์เช่นเดียวกับ Zeus แต่แตกต่างกันที่เส้นทาง (Path) ที่ร้องขอการเชื่อมต่อเนื่องจากมีการตั้งค่าไม่เหมือนกัน นอกจากนี้ขนาดข้อมูลที่มีการสื่อสารกันและคาบเวลาที่ใช้ในการสื่อสารยังต่างกัน โดยที่ VertexNet มีการสร้างการเชื่อมต่อทุกๆ 30 วินาทีและมีขนาดข้อมูลคือ 193 ไบต์ดังรูป

No.	Time	Source	Destination	Protocol	Source Port	Destination Port	Host	Total Length	Info
27028	42.27020	161.246.5.184	161.246.5.242	HTTP	1419	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
37018	42.43910	161.246.5.184	161.246.5.242	HTTP	2413	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
47008	42.5887	161.246.5.184	161.246.5.242	HTTP	3402	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
57001	42.75740	161.246.5.184	161.246.5.242	HTTP	4411	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
67010	42.92610	161.246.5.184	161.246.5.242	HTTP	5423	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
77015	43.09480	161.246.5.184	161.246.5.242	HTTP	6438	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
87019	43.26350	161.246.5.184	161.246.5.242	HTTP	7453	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
97028	43.43220	161.246.5.184	161.246.5.242	HTTP	8468	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
10737	43.60090	161.246.5.184	161.246.5.242	HTTP	9483	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
11746	43.76960	161.246.5.184	161.246.5.242	HTTP	10498	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
12755	43.93830	161.246.5.184	161.246.5.242	HTTP	11513	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
13764	44.10700	161.246.5.184	161.246.5.242	HTTP	12528	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
14773	44.27570	161.246.5.184	161.246.5.242	HTTP	13543	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
15782	44.44440	161.246.5.184	161.246.5.242	HTTP	14558	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
16791	44.61310	161.246.5.184	161.246.5.242	HTTP	15573	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
17800	44.78180	161.246.5.184	161.246.5.242	HTTP	16588	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
18809	44.95050	161.246.5.184	161.246.5.242	HTTP	17603	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
19818	45.11920	161.246.5.184	161.246.5.242	HTTP	18618	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
20827	45.28790	161.246.5.184	161.246.5.242	HTTP	19633	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
21836	45.45660	161.246.5.184	161.246.5.242	HTTP	20648	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
22845	45.62530	161.246.5.184	161.246.5.242	HTTP	21663	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
23854	45.79400	161.246.5.184	161.246.5.242	HTTP	22678	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
24863	45.96270	161.246.5.184	161.246.5.242	HTTP	23693	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
25872	46.13140	161.246.5.184	161.246.5.242	HTTP	24708	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
26881	46.30010	161.246.5.184	161.246.5.242	HTTP	25723	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
27890	46.46880	161.246.5.184	161.246.5.242	HTTP	26738	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
28899	46.63750	161.246.5.184	161.246.5.242	HTTP	27753	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
29908	46.80620	161.246.5.184	161.246.5.242	HTTP	28768	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
30917	46.97490	161.246.5.184	161.246.5.242	HTTP	29783	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
31926	47.14360	161.246.5.184	161.246.5.242	HTTP	30798	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
32935	47.31230	161.246.5.184	161.246.5.242	HTTP	31813	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
33944	47.48100	161.246.5.184	161.246.5.242	HTTP	32828	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
34953	47.64970	161.246.5.184	161.246.5.242	HTTP	33843	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
35962	47.81840	161.246.5.184	161.246.5.242	HTTP	34858	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
36971	47.98710	161.246.5.184	161.246.5.242	HTTP	35873	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
37980	48.15580	161.246.5.184	161.246.5.242	HTTP	36888	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
38989	48.32450	161.246.5.184	161.246.5.242	HTTP	37903	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
39998	48.49320	161.246.5.184	161.246.5.242	HTTP	38918	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
41007	48.66190	161.246.5.184	161.246.5.242	HTTP	39933	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
42016	48.83060	161.246.5.184	161.246.5.242	HTTP	40948	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
43025	49.00000	161.246.5.184	161.246.5.242	HTTP	41963	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
44034	49.16870	161.246.5.184	161.246.5.242	HTTP	42978	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
45043	49.33740	161.246.5.184	161.246.5.242	HTTP	43993	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
46052	49.50610	161.246.5.184	161.246.5.242	HTTP	45008	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
47061	49.67480	161.246.5.184	161.246.5.242	HTTP	46023	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
48070	49.84350	161.246.5.184	161.246.5.242	HTTP	47038	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
49079	50.01220	161.246.5.184	161.246.5.242	HTTP	48053	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
50088	50.18090	161.246.5.184	161.246.5.242	HTTP	49068	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
51097	50.34960	161.246.5.184	161.246.5.242	HTTP	50083	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
52106	50.51830	161.246.5.184	161.246.5.242	HTTP	51098	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
53115	50.68700	161.246.5.184	161.246.5.242	HTTP	52113	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
54124	50.85570	161.246.5.184	161.246.5.242	HTTP	53128	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
55133	51.02440	161.246.5.184	161.246.5.242	HTTP	54143	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
56142	51.19310	161.246.5.184	161.246.5.242	HTTP	55158	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
57151	51.36180	161.246.5.184	161.246.5.242	HTTP	56173	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
58160	51.53050	161.246.5.184	161.246.5.242	HTTP	57188	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
59169	51.69920	161.246.5.184	161.246.5.242	HTTP	58203	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
60178	51.86790	161.246.5.184	161.246.5.242	HTTP	59218	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
61187	52.03660	161.246.5.184	161.246.5.242	HTTP	60233	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
62196	52.20530	161.246.5.184	161.246.5.242	HTTP	61248	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
63205	52.37400	161.246.5.184	161.246.5.242	HTTP	62263	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
64214	52.54270	161.246.5.184	161.246.5.242	HTTP	63278	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
65223	52.71140	161.246.5.184	161.246.5.242	HTTP	64293	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
66232	52.88010	161.246.5.184	161.246.5.242	HTTP	65308	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
67241	53.04880	161.246.5.184	161.246.5.242	HTTP	66323	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
68250	53.21750	161.246.5.184	161.246.5.242	HTTP	67338	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
69259	53.38620	161.246.5.184	161.246.5.242	HTTP	68353	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
70268	53.55490	161.246.5.184	161.246.5.242	HTTP	69368	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
71277	53.72360	161.246.5.184	161.246.5.242	HTTP	70383	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
72286	53.89230	161.246.5.184	161.246.5.242	HTTP	71398	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
73295	54.06100	161.246.5.184	161.246.5.242	HTTP	72413	80	161.246.5.242	193	GET /1/cfg.bin HTTP/1.1
74304	54.22970	161.246.5.184	161.246.5.242	HTTP	73428</				

อย่างไรก็ตามจากการวิเคราะห์แพ็คเก็ตทั้งหมดทำให้สามารถสรุปได้ว่าการเชื่อมต่อของบอท จะมีการปะปนอยู่กับการเชื่อมต่อของเครื่องปกติแต่จะมีการเชื่อมต่อบางอย่างของบอทที่จะมีลักษณะที่จำเพาะและมีความซับซ้อนที่สามารถตรวจสอบได้โดยการพัฒนาโปรแกรมเพื่อตรวจจับพฤติกรรมเหล่านี้

#### 4.5 การทดลองการคำนวณตามเอกสารอ้างอิง [1]

ในการวิเคราะห์พฤติกรรมของเครือข่ายนั้นจำเป็นต้องมีการคำนวณพารามิเตอร์ที่ต่างชนิดกันซึ่งสามารถออกได้เป็น 3 ชนิดคือ ข้อมูลเชิงเลขคณิต เช่น จำนวนไบต์ ข้อมูลเชิงประเภท เช่น แฟล็กของทีซีพีเฮดเดอร์ และข้อมูลเชิงลำดับชั้น เช่น หมายเลขไอพี เป็นต้น ในการคำนวณและเปรียบเทียบข้อมูลแต่ละชนิดจะมีกระบวนการที่แตกต่างกันซึ่งได้กล่าวไว้ในบทที่สาม

ในการทดลองการคำนวณตามอัลกอริทึมเบิร์ชมันน์ได้กำหนดข้อมูล 3 แพ็คเก็ตซึ่งเป็นข้อมูลที่ได้มีการเก็บจากสถานการณ์ปกติ 1 แพ็คเก็ตและเก็บจากแพ็คเก็ตของ VertexNet อีก 2 แพ็คเก็ตซึ่งข้อมูลแต่ละชุดมีรูปแบบดังนี้

Source Address, Destination Address, Source Port, Destination Port, Bytes, Protocol

1. 161.246.5.183, 161.246.5.182, 1050, 80, 40 TCP
2. 161.246.5.184, 161.246.5.182, 1055, 80, 144 UDP
3. 161.246.5.86, 38.111.126.203, 4804, 23468, 144 TCP

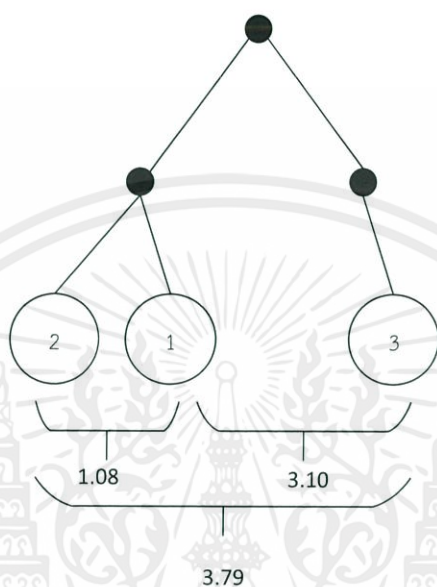
กำหนดให้ เริ่มต้นทำการพิจารณาแพ็คเก็ตที่ 1 ทำการสร้างเป็น CF Tree ที่มี 1 โหนด จากนั้นทำการเพิ่มแพ็คเก็ตที่ 2 เข้ามาและสร้างเป็น CF Tree ที่มี 2 โหนด จากนั้นจึงพิจารณาโหนดที่ 3 กำหนด  $N = 2$  หมายความว่า 1 โหนดจะมีสมาชิกได้ไม่เกิน 2 โหนดเท่านั้น กำหนด  $TCP = \{0,1\}$ ,  $UDP = \{1,0\}$  เมื่อทำการคำนวณ CF Entry ของแต่ละชุดข้อมูลและจะสามารถเริ่มเปรียบเทียบแต่ละ Entry เพื่อสร้าง CF Tree ได้

N	Source Address		Destination Address		SrcPort		DstPort		Bytes	Protocol
	c	Radius	c	Radius	c	Range	c	Range		
1	161.246.5.183/32	161.246.5.183	161.246.5.182/32	161.246.5.182	1050/16	1050	80/16	80	40	{01}
1	161.246.5.183/32	161.246.5.183	161.246.5.184/32	161.246.5.184	1055/16	1055	80/16	80	144	{10}
1	38.111.126.203/32	38.111.126.203	161.246.5.86/32	161.246.5.86	4804/16	4804	23168/16	23468	144	{01}

รูปที่ 4.13 ผลการคำนวณ CF Entry ของข้อมูลแต่ละชุด

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากการคำนวณสามารถหาค่าความแตกต่างระหว่างชุดข้อมูลได้ซึ่งผลที่ได้คือ ข้อมูลชุดที่ 1 มีความใกล้เคียงกับชุดที่ 2 มากกว่าชุดที่ 3 และข้อมูลชุดที่ 3 มีความแตกต่างจากชุดที่ 2 มากกว่าชุดที่ 1 ดังนั้นทำให้สามารถจัดเป็น CF Tree ได้ดังนี้



รูปที่ 4.14 การจำลอง CF Tree ที่ได้จากการคำนวณ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## บทที่ 5

# บทสรุปและข้อเสนอแนะ

จากการศึกษาการแนวทางการทำงานต่างๆเพื่อหาวิธีในการตรวจจับบอทเน็ตอย่างมีประสิทธิภาพทำให้ทราบถึงกระบวนการต่างๆมากมายที่ถูกประยุกต์ใช้เพื่อจุดประสงค์ดังกล่าว สำหรับการตรวจจับบอทเน็ตโดยใช้กระบวนการสังเกตพฤติกรรมการใช้งานเครือข่ายของคอมพิวเตอร์ต่างๆภายในเครือข่ายนั้นมีปัญหาหลายอย่างซึ่งทำให้เกิดปัญหาของการศึกษาค้นคว้าเป็นต้นว่า การแบ่งกลุ่มข้อมูล การทำงานกับข้อมูลที่มีจำนวนมากหรือแม้แต่การจัดการความเร็วในการประมวลผล เป็นต้น

### 5.1 บทสรุป

โครงการเรื่องระบบตรวจจับบอทเน็ตจากพฤติกรรมการใช้งานเครือข่ายได้ถูกออกแบบมาเพื่อใช้ในการตรวจจับบอทเน็ตจากการสังเกตการใช้งานเครือข่ายที่ใกล้เคียงกับกลุ่มข้อมูลตัวอย่างซึ่งแตกต่างจากเทคนิคการตรวจจับโดยใช้ลักษณะเฉพาะของบอทเน็ตเป็นสิ่งบ่งชี้ถึงความน่าจะเป็นที่จะเป็นบอท ข้อดีของการใช้การสังเกตพฤติกรรมแทนการใช้ลักษณะบ่งชี้ดังกล่าวคือความสามารถในการตรวจหาบอทจากพฤติกรรมใกล้เคียงโดยจะสามารถตรวจจับบอทเน็ตได้ถึงแม้ว่าจะไม่มีการอัปเดตฐานข้อมูลของกลุ่มตัวอย่างเป็นประจำอย่างการใช้งานแอนตี้ไวรัสหรือระบบป้องกันเครือข่ายๆแบบอื่นๆ ซึ่งจะทำให้สามารถป้องกันความเสี่ยงที่อาจเกิดขึ้นได้อย่างกว้างขวางมากขึ้นกว่าเดิม

### 5.2 ปัญหา อุปสรรค และแนวทางการแก้ไข

- 1) เอกสารที่เกี่ยวกับการทำงานของอัลกอริทึมเบียร์ชิ่งมีจำนวนน้อยและส่วนมากเป็นชิ้นงานนำเสนอซึ่งคัดลอกกันมาเป็นทอดๆ ทำให้สามารถศึกษากรณีตัวอย่างและทำความเข้าใจได้ยากมากขึ้นเพราะมีจำนวนกรณีตัวอย่างที่แตกต่างกันน้อย
- 2) การศึกษาพฤติกรรมของการทำงานของเครือข่ายจะได้ข้อมูลซึ่งส่วนมากไม่ใช่ข้อมูลที่สามารถนำไปคำนวณทางคณิตศาสตร์ได้โดยตรง ผู้วิจัยจึงจำเป็นต้องมีการศึกษาการแปลงหรือแทนค่าและกำหนดกระบวนการทางคณิตศาสตร์ให้กับข้อมูลดังกล่าวอย่างเหมาะสมก่อนที่จะนำไปคำนวณ

เอกสารนี้เป็นเอกสารต้นฉบับไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณี (3) เป้าหมายของโครงการคือระบบที่สามารถตรวจจับบอทเน็ตได้อย่างมีประสิทธิภาพโดยใช้การสังเกตพฤติกรรมและการแบ่งกลุ่มข้อมูลของการทำงานของเครือข่ายดังนั้นจึงต้องมีการทำงานกับ

ข้อมูลจำนวนมากและใช้เวลานานในการเฝ้าสังเกตพฤติกรรมทำให้ระบบไม่เป็นแบบตามเวลาจริง

- 4) การวิเคราะห์การใช้งานเครือข่ายของโครงงานนี้อาศัยการทำงานของโปรโตคอล เน็ตโพล์ ที่พัฒนาโดยซิสโก้แต่เนื่องจากอุปกรณ์ของสาขาวิชาไม่รองรับการใช้งานดังกล่าว ผู้วิจัยจึงต้องมีการศึกษาการใช้ซอฟต์แวร์อื่นทดแทนที่สามารถสร้างแพ็คเก็ต เน็ตโพล์ เพื่อให้ข้อมูลที่ระบบต้องการได้
- 5) การบันทึกพฤติกรรมบอทเน็ตจำเป็นต้องทำอย่างยิ่งที่ต้องมีการศึกษาพฤติกรรมในระบบปิดก่อนซึ่งต้องมีการติดตั้งและศึกษาจากเวอร์ชวลแมชีนหลายเครื่องซึ่งใช้ทรัพยากรระบบเป็นอย่างมาก
- 6) การศึกษาบอทเน็ตนั้นต้องใช้บอทเน็ตจริงหรือใกล้เคียงซึ่งหาได้ยากและง่ายต่อการถูกหลอกลงให้ติดไวรัสหรือมัลแวร์

### 5.3 แนวทางในการพัฒนาต่อ

- 1) นำระบบที่พัฒนาได้ไปผนวกกับระบบเรียนรู้ของคอมพิวเตอร์ซึ่งจะทำให้ระบบสามารถตรวจจับบอทเน็ตหรือมัลแวร์อย่างอื่นได้อย่างมีประสิทธิภาพ
- 2) สามารถนำระบบไปพัฒนาเพื่อวิเคราะห์และแบ่งกลุ่มการใช้งานเครือข่ายเพื่อตรวจหาความผิดปกติของการใช้งานเครือข่ายได้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

## บรรณานุกรม

- [1] Abdun Naser Mahmood, Christopher Leckie and Parampali Udaya. “An Efficient Clustering Scheme to Exploit Hierarchical Data in Network Traffic Analysis” IEEE Trans. On Knowle, vol. 20, no. 6, 2007
- [2] T. Zhang, R. Ramakrishnan and M. Livny. “Birch: An Efficient Data Clustering Method for Very Large Databases” ACM, inc. 1996. pp103-114
- [3] Meisam E. “botAnalytics: Improving HTTP-Based Botnet Detection By Using Network Behavior Analysis System” M. Thesis of Faculty of Computer Science and Information Technology of University of Malaya. 2010.
- [4] Alexander B. “Network Characterization For Botnet Detection Using Statistical-Behavioral Methods” M. Thesis of Thayer School of Engineering Dartmouth College Hanover. 2009
- [5] Yuanyuan Z. “On Detection of Current and Next-Generation Botnets” Ph.D Thesis of Michigan of University. 2012
- [6] Jae-Seo, HyunCheol. “The Activity of Malicious HTTP-based Botnets using Degree of Periodic Repeatability” 83-86. IEEE Trans. 2008

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า  
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้