# Euclidian- and Cosine-Distances based Detection of Distributed Host Search Attacks

Yasuo Musashi,[†] Satoshi Dobayashi[††] Dennis Arturo Ludeña Romaña,[†*] Shinichiro Kubota,[†] and Kenichi Sugitani[†]

We statistically investigated the total PTR resource record (RR) based DNS query request packet traffic from the Internet to the top domain DNS server in a university campus network through January 1st to December 31st, 2011.   The obtained results are: (1) We found twelve host search (HS) attacks in the scores for detection method using the calculated Euclidean distances between the observed IP address and the last observed IP address in the DNS query keywords by employing both threshold ranges of 1.0-2.0 (consecutive) and 150.2-210.4 (random).   However, we found nineteen HS attacks in the scores using the calculated cosine distance between the DNS query IP addresses (threshold ranges of 0.75-0.83 and 0.9-1.0).   (3) In the newly found HS attacks, we observed that the source IP addresses of the HS attack DNS query packets are distributed Therefore, it can be concluded that the cosine distance based detection technology can detect the source IP address-distributed host search attack.

## 1.  Introduciton

   It is of considerable importance to raise up a detection rate of bots, since they become components of the bot clustered networks that are used to transmit a lot of unsolicited mails including like spam, phishing, and spam mailing activities and to execute distributed denial of service attacks [1-4].

   The HS attack is recognized to be a pre-investigation activity or a harvesting attack of fully qualified domain names (FQDNs) of the university campus and/or enterprise networks *i.e.* after the HS attack, the attacker can concentrate to check out the vulnerabilities in the targeted servers or hosts.

   Previously, we reported development and evaluation of the Euclidian distance based detection model system for the HS attack against the campus top domain name system (DNS) server [5] and it currently works still well for detecting the single source IP address-based HS attack, however, recently, the HS attackers upgraded their strategies that they started to a distributed source IP address HS attack to evade the detection system *i.e.* it is required to develop a new system for detecting the distributed source IP address based HS attack.
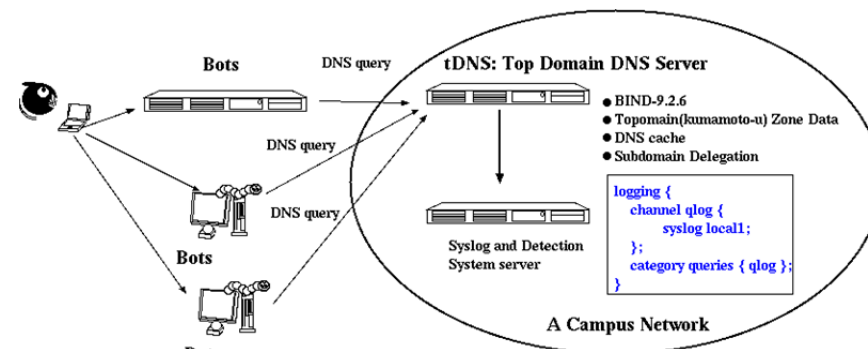
Figure 1   A schematic diagram of a network observed in the present study.

   In this paper, (1) we carried out Euclidean and cosine distances based analyses on the source IP address and the DNS query IP address in the total PTR resource record (RR) based DNS query request packet traffic from the Internet through January 1st to December 31st, 2011, and (2) we assessed the both results for the Euclidean and cosine distances based analyses on the IP addresses as the query keywords in the PTR-RR based DNS query packet traffic.

## 2.  Observation

### 2.1  Network Systems and DNS Query Packet Capturing

   We investigated on the DNS query request packet traffic between the top domain (tDNS) DNS server and the DNS clients.   Figure 1 shows an observed network system in the present study, which consists of the tDNS server and the PC clients as bots like a host search bot or a spam bot in the campus or enterprise network, and the victim hosts like the DNS servers on the campus network.   The tDNS server is one of the top level domain name (kumamoto-u) system servers and plays an important role of domain name resolution including DNS cache function, and subdomain name delegation services for many PC clients and the subdomain network servers, respectively, and the operating system is Linux OS (CentOS 5.5 Final) in which the kernel-2.6.18 is currently employed with the Intel Xeon X5660 2.8 GHz 6 Cores dual node system, the 16GB core memory, and Intel Corporation EthernetPro 82575EB Gigabit Ethernet Controller.

   In the tDNS server, the BIND-9.3.6-P1 program package has been employed as a DNS server daemon [6].   The DNS query request packet and their query keywords have been captured and decoded by a query logging option (see Figure 1 and the named.conf manual of the BIND program in more detail).   The log of DNS query request packet access has been recorded in the syslog files.   All of the syslog files are daily updated  by  the  cron  system.
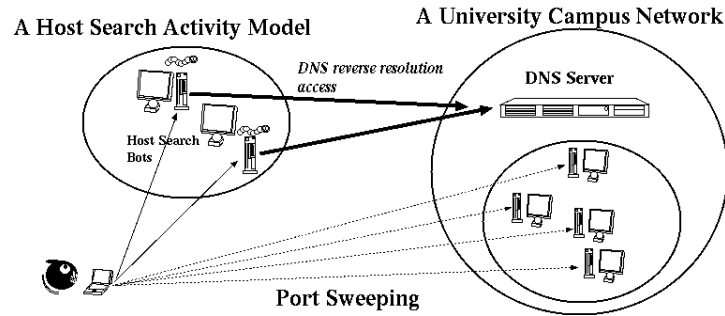
Figure 2    A host search (HS) attack model.

The line of syslog message consists of the contents of the DNS query request packet like a time, a source IP address of the DNS client, a fully qualified domain name (A and AAAA resource record (RR) for IPv4 and IPv6 addresses, respectively) type, an IP address (PTR RR) type, or an E-mail exchange (MX RR) type.

**2.2    Host Search Attack Model**

We define here a host search (HS) model (See Figure 2).
— *A host search (HS) attack model* — the host search (HS) attack can be mainly carried out by a small number of IP hosts on the Internet or in the campus network like bot compromised PCs or the public cloud.    Since these IP hosts send a lot of the DNS reverse name resolution (the PTR RR based DNS query) request packets to the tDNS server, the unique IP addresses- and the unique DNS query-keywords based entropies decrease and increase, simultaneously [7].

Here, we should also define thresholds for detecting the HS attack, as setting to 10,000 packets day$^{-1}$ for the frequencies of the top ten unique source IP addresses or the DNS query keywords.

We also investigated the IP address change in the PTR RR based DNS query request packet traffic through January 8th and 21st, 2009, and the results are shown in Figure 3.    In Figure 3A, at January 8th, 2009, we can view scenery that the IP address as DNS query keyword is consecutively incremented.    Therefore, it has a possibility that this consecutive increment of the IP address can be useful to detect the HS attack in the PTR RR based DNS query request packet traffic (consecutive model).    In Figure 3B, at January 21st, 2009, we can see it that the IP address as DNS query keyword is discontinuously or randomly changed (random model).

From these results, we need to take into consideration on the consecutive and the random IP address query keyword based models in order to develop an HS attack detection system.

**2.3    Euclidean- and Cosine-Distances of IP addresses as DNS Query Keywords**

The Euclidean distances, **qd(IP$_i$, IP$_{i-1}$)**, are calculated, as
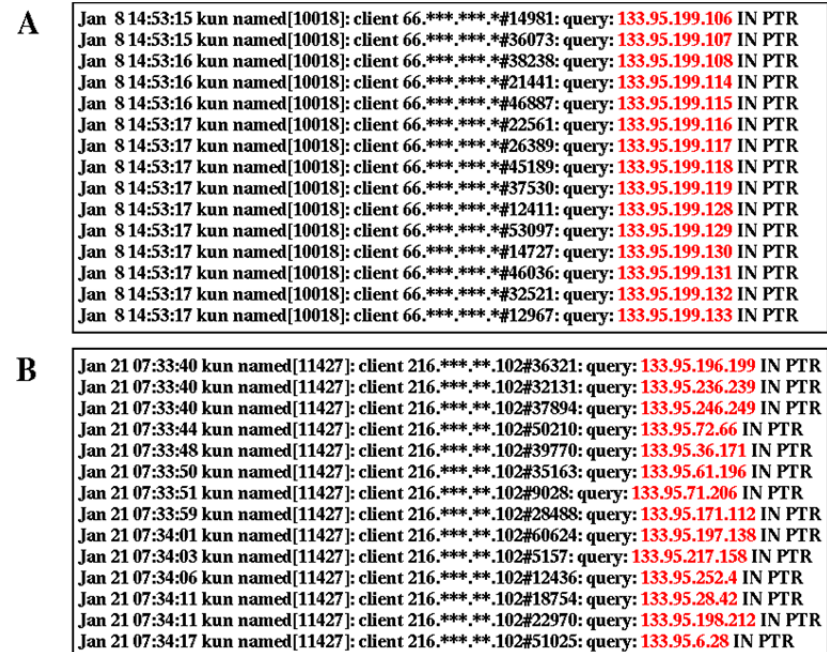


Figure 3    Changes in the IP address as the DNS query keywords in the total PTR-resource records (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server at January 8th (A) and 21st (B), 2009.

$$d(IP_i, IP_{i-1}) = \sqrt{\sum_{j=1}^{4} (x_{i,j} - x_{i-1,j})^2} \qquad (1)$$

where both IP$_i$ and IP$_{i-1}$ are the current IP address i and the last IP address i-1 of the DNS query keywords, respectively, and where x$_{i,1}$, x$_{i,2}$, x$_{i,3}$, and x$_{i,4}$ correspond to an IPv4 address like A.B.C.D, respectively.    For instance, if an IP address is 192.168.1.1, the vector (x$_{i,1}$, x$_{i,2}$, x$_{i,3}$, x$_{i,4}$) can be represented as (192.0, 168.0, 1.0, 1.0).

If the HS attack model follows the consecutive DNS query keyword based model, the detection is decided by thresholds qd$_{min}$=1.0 and qd$_{max}$=5.0 [8], as

$$qd_{min} (= 1.0) \leq qd(IP_i, IP_{i-1}) \leq qd_{max} (= 5.0) \qquad (2)$$

The campus IP addresses are represented as 133.95.x$_i$.y$_i$ in which both x$_i$ and y$_i$ can take numbers from 0 to 255, as: $0 \leq x_i \leq 255$ and $0 \leq y_i \leq 255$, and $(x_i - x_{i-1})^2$ or $(y_i - y_{i-1})^2$ takes a
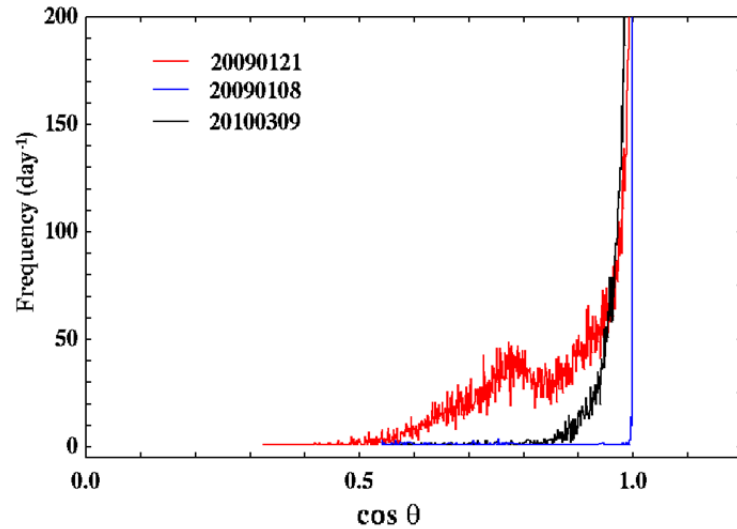
Figure 4    Frequency distributions of the cosine distance at January 10th, 21st, and March 9th, 2010 (day$^{-1}$ unit).

range from 0 to 255$^2$ *i.e.* the range of the Euclidian distance, **qd(IP$_i$, IP$_{i-1}$)**, should be from 0.0 to $\sqrt{255^2 + 255^2}$ (~ 360.6). If the Euclidian distance, **qd(IP$_i$, IP$_{i-1}$)**, follows the Gaussian distribution, the probability for the Euclidian distance takes a maximum value between at 180.3 (~ 360.6/2) with a standard deviation of 30.1 (~ 360.6/12) because of the central limit theorem *i.e.* qd$_{min}$ and qd$_{max}$ should take values of 150.2 (~ 180.3-31.1) and 210.4 (~ 180.3+31.1) [9,10].

$$qd_{min} (= 150.2) \leq qd(IP_i, IP_{i-1}) \leq qd_{max} (= 210.4) \qquad (3)$$

The cosine distances, **cos θ (IP$_i$, IP$_{i-1}$)**, are obtained, as

$$\cos \theta(IP_i, IP_{i-1}) = \frac{IP_{i-1}^T \bullet IP_i}{\left| IP_{i-1}^T \right| \left| IP_i \right|} \qquad (4)$$

where the cosine distance takes a range from 0.413 (=**cos θ $_{min}$**) to 1.000, since **cos θ $_{min}$** is estimated from cosine distance between vectors $(133, 95, 0, 0)^T$ and $(13.95.255.255)^T$.

In Figure 4, we show the calculated frequency distribution of the cosine distances at January 8th (consecutive model), 21st (random model; normal distribution), 2009, and March 9th,

```
1   #!/bin/tcsh -f
2   set Threshold=10
3   # Step 1 Reduction of the Noise
4   cat /var/log/querylog | clgrep -v -cclients.conf  | \
5   grep "IN PTR" | arpa | \
6   awk '{print $9}' | sort -r | uniq -c | sort -r | \
7   awk '{printf("%s\t%s\n",$2,$1);}' | \
8   qdos 1000 >noise.conf
9   # Step 2 Learning to produce a low-diemnsianl
10  cat /var/log/querylog | clgrep -v -cclients.conf | \
11  grep "IN PTR" | arpa | \
12  cngrep -v -Dnoise.conf | \
13  sdis 0.0 0.0 | \
14  qdis 1.0 5.0 150.2 210.4 | \
15  tr '#' ' ' | awk '{print $7}' | sort -r | uniq -c | sort -r | \
16  awk '{printf("%s\t%s\n",$2,$1);}' | qdos $Threshold | \
17  awk '{print $1}' >tmpfile
18  # Step 3 Detection
19  cat /var/log/querylog | clgrep -ctmpfile | \
20  grep "IN PTR" | arpa >HSdet.log
21  # Step 4 Scoring
22  cat HSdet.log | wc -l >>HSdetScore.txt
23  exit 0
```

Figure 5    Host Search Attack Detection Algorithm for Euclidian distance.

2010 (random model; exponential distribution). The frequency distribution at January 21st, 2009, has a significant peak at a range between 0.73 and 0.83 as well as quick increasing from 0.9 to 1.0, indicating that the thresholds cd$_{min}$ and cd$_{max}$ should take ranges between 0.73-0.83 and 0.9-1.0.

$$cd_{min} (= 0.73 \, or \, 0.9) \leq \cos \theta (IP_i, IP_{i-1}) \leq cd_{max} (= 0.9 \, or \, 1.0) \qquad (5)$$

**2.4  Detection Algorithm for Host Search Activity**

We suggest the following detection algorithm of the Host Search (HS) activity and we show

```
13  sdis 0.0 0.0 | \
14  qdis -yz 0.73 0.83 0.9 1.0 | \
```

Figure 6    Host Search Attack Detection Algorithm for Cosine Distance.

a prototype program (see Figure 5):

⎯⎯ **Step 1**    *Reduction of the Noise*⎯⎯ In this step, the **clgrep** and **grep** commands extract the inbound PTR RR based DNS query request packet messages from the DNS query log file (*/var/log/querylog*), the **arpa** command converts the reverse query format "D.C.B.A.in-addr.arpa" into the usual IPv4 format "A.B.C.D" (A, B, C, and D represent digit numbers of {0-255}), the top **awk** commands and the two **sort** and one **uniq** commands calculate and print out the IP address query keywords and their frequencies, and the **qdos** command prints out the IP addresses and the frequencies into the *noise.conf* file when the frequencies are greater than 1,000 day$^{-1}$.

⎯⎯ **Step 2**    *Learning to produce a low-dimensional*⎯⎯ In this step, t the **clgrep**, **grep, arpa,** commands take the same functions as ones in **Step 1**, the **cngrep** command discards the IP addresses listed in the *noise.conf* file from the syslog messages, the **sdis** command prints out a syslog message if the Euclidean distance **sd(IP$_i$, IP$_{i-1}$)** between the two source IP addresses is calculated to be zero, the **qdis** command prints out the syslog message if the Euclidean distance **qd(IP$_i$, IP$_{i-1}$)** takes ranges of 1.0-2.0 and 150.2-210.4, or the **qdis -yz** command (Figure 6) pints out the syslog messages if the cosine distance **cos θ (IP$_i$, IP$_{i-1}$)** takes ranges of 0.73-0.82 and 0.9-1.0, and the **awk**, **sort**, **uniq**, and **qdos** commands (lines 15 to 17 in Figure 5) compute the frequencies of the Euclidean distance **qd(IP$_i$, IP$_{i-1}$)** and if the frequency exceeds a threshold value (*Threshold*=10), they write out the candidate IP addresses into a *tmpfile* as training data.

⎯⎯ **Step 3** *Detection* ⎯⎯ In the next step, the **clgrep**, **grep**, and **arpa** commands extract the HS activity related messages in the DNS query log file (*/var/log/querylog*), using the training data (*tmpfile*) and they generate only an HS activity related DNS query log file (*HSdet.log*).

⎯⎯ **Step 4** *Scoring* ⎯⎯In the final step, the **wc** command calculates the score for the detection of the HS activity in the file *HSdet.log*, and it writes out the detection score into a score file (*HSdetScore.txt*).

## 3.   Results and Discussion

### 3.1   Euclidean distance- and cosine distance-based Host Search Detection Model

We illustrate the calculated score of the host search (HS)   attack  using  Euclidean-  and
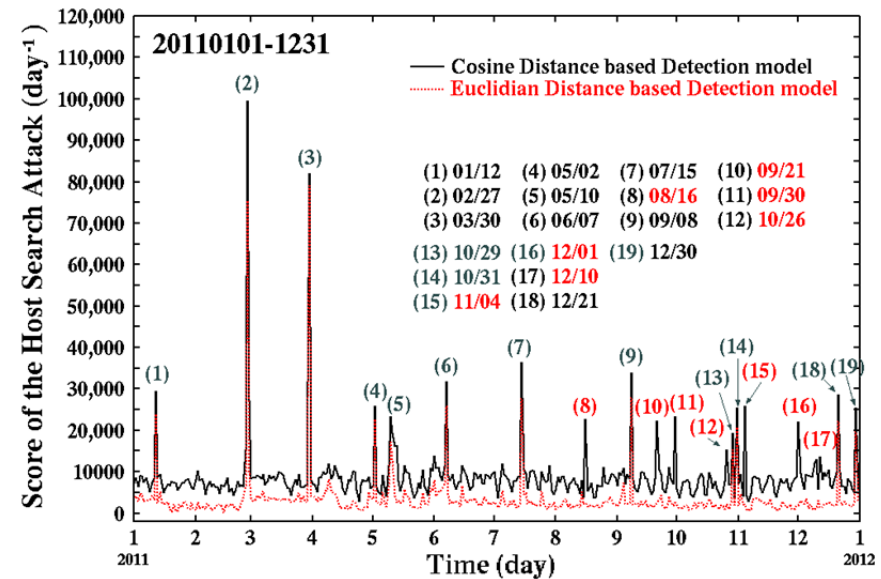


Figure 7    Changes in score of the host search (HS) attack detection in the total PTR resource records (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server through January 1st to December 31st, 2011 (day$^{-1}$ unit).   The solid- and the dotted-curves show scores for the cosine distance- and the Euclidian distance-based detection models, respectively.

cosine-distance based detection models (1.0 ≤ **d(IP$_i$, IP$_{i-1}$)** ≤ 2.0 or 150.2 ≤ **d(IP$_i$, IP$_{i-1}$)** ≤ 210.4.) or ((0.73 ≤ **cos θ (IP$_i$, IP$_{i-1}$)** ≤ 0.83 or 0.9 ≤ **cos θ (IP$_i$, IP$_{i-1}$)** ≤ 1.0)) between the current IP address IP$_i$ and the last IP address IP$_{i-1}$, as the DNS query keywords in the PTR resource record (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server through January 1st to December 31st, 2011, as shown in Figure 7.

In Figure 7, we can observe nineteen significant peaks (1)-(19) being allocated to (1) January 12th, (2) February27th, (3) March 30th, (4) May 2nd, (5) 10th, (6) June 7th, (7) July 15th, (8) August 16th, (9) September 8th, (10) 21st, (11) 30th, (12) October 26th, (13) 29th, (14) 31st, (15) November 4th, (16) December 1st, (17) 10th, (18) 21st, and (19) 30th, 2011, respectively.

In the score curve for the Euclidian distance based detection model, we can find no peaks corresponding to the peaks (8), (10), (11), (12), (15), (16), and (17) in the score curve for the cosine distance based detection model.   This result shows   that   the   cosine   distance   based

```
Aug 16 01:19:24 kun named[13868]: client ***.125.92.90#64965: query: 133.95.25.19 IN PTR
Aug 16 01:19:28 kun named[13868]: client ***.125.92.88#46342: query: 133.95.25.25 IN PTR
Aug 16 01:19:32 kun named[13868]: client ***.125.90.83#57923: query: 133.95.25.31 IN PTR
Aug 16 01:19:32 kun named[13868]: client ***.125.90.82#54527: query: 133.95.25.32 IN PTR
Aug 16 01:19:33 kun named[13868]: client ***.125.90.91#44176: query: 133.95.25.34 IN PTR
Aug 16 01:19:33 kun named[13868]: client ***.125.92.83#59507: query: 133.95.25.35 IN PTR
Aug 16 01:19:33 kun named[13868]: client ***.125.92.82#64224: query: 133.95.25.38 IN PTR
Aug 16 01:19:34 kun named[13868]: client ***.125.90.86#46685: query: 133.95.25.39 IN PTR
Aug 16 01:19:34 kun named[13868]: client ***.125.90.88#53614: query: 133.95.25.40 IN PTR
Aug 16 01:19:35 kun named[13868]: client ***.125.92.90#45848: query: 133.95.25.43 IN PTR
Aug 16 01:19:36 kun named[13868]: client ***.125.90.82#57662: query: 133.95.25.45 IN PTR
Aug 16 01:19:37 kun named[13868]: client ***.125.90.87#53354: query: 133.95.25.44 IN PTR
Aug 16 01:19:40 kun named[13868]: client ***.125.92.91#53858: query: 133.95.25.50 IN PTR
Aug 16 01:19:40 kun named[13868]: client ***.125.92.83#49809: query: 133.95.25.49 IN PTR
Aug 16 01:19:41 kun named[13868]: client ***.125.92.85#48131: query: 133.95.25.51 IN PTR
```

Figure 8   Changes in the source IP address in the total PTR-resource records (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server at August 16th, 2011.
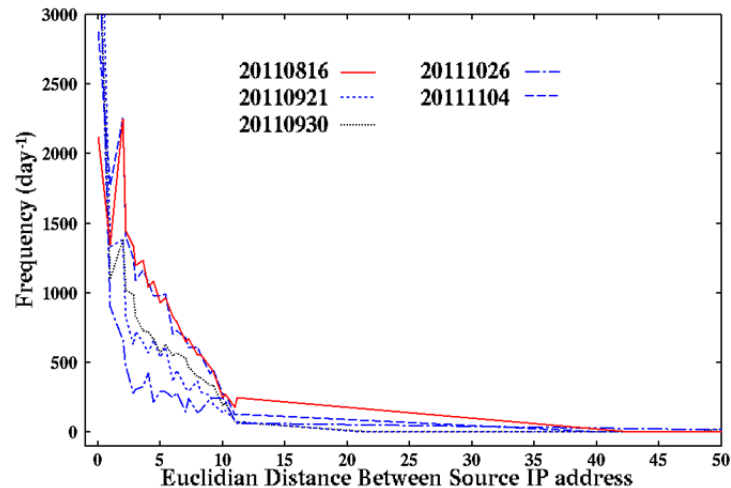


Figure 9   Frequency distributions of the Euclidian distance between the source IP addresses at August 16th, September 21st, 30th, October 26th, and November 4th, 2011 (day$^{-1}$ unit).

detection technology is much precisely than the Euclidian distance based one or much false positive.   Thus, we also investigated the source IP address- and query IP address-changes in the PTR RR based DNS query request packet traffic through August 16th, 2011, and the

results are shown in Figure 8.

In Figure 8, we can view scenery that the source IP addresses change periodically and the query IP addresses are incremented like a consecutive manner, showing that the cosine distance based detection model can be useful for detecting the source IP address distributed host search (HS) attack like a distributed denial of service (DDoS) attack.

**3.2   Frequency Distribution of the Euclidian distance in Source IP addresses**

We calculated frequency distributions of the Euclidian distance for the five peaks (8), (10), (11), (12), (15), as shown in Figure 9.   In Figure 9, the each frequency distribution has a significant peak and all the peaks take a range of 1.0-5.0.   The detection of the source IP address distributed attack is decided by thresholds of $sd_{min}=1.0$ and $sd_{max}=5.0$, as

$$sd_{min}(=1.0) \le sd(IP_i, IP_{i-1}) \le sd_{max}(=5.0) \qquad (6)$$

```
13  sdis 0.0 0.0 1.0 5.0 | \
14  qdis 1.0 5.0 150.2 210.4 | \
```

Figure 9   Source IP address Distributed Host Search Attack Detection Algorithm.

addresses as the DNS query keywords, it can improve the Euclidian distance based HS attack detection technology *i.e.* we can raise the HS detection rate but decrease the false positive.

From these results, we show the newly improved Euclidian distance based HS attack detection technology in the next section of 3.3.

**3.3   Improved Detection Algorithm for Source IP address Distributed Host Search Attack**

We suggest again the following detection algorithm of the source IP address distributed host search (HS) attack and we show a new prototype program (see Figures 5 and 9):

── **Step 1** ──This step is as the completely same as **Step 1** in Figure 5.

── **Step 2** ── In this step, the **sdis** command is only different from that one in **Step 1** in Figure 9 (See also Figure 5), in which **sdis** command prints out the syslog message if the source IP address Euclidean distance **sd(IP$_i$, IP$_{i-1}$)** takes zero or a range of 0. and 1.0-5.0.

── **Steps 3 and 4** ──These two steps are as the completely same as **Steps 2** and **3** in Figure 5, respectively.

**3.4   Evaluation**

We calculated the score for the newly improved HS attack detection model in the inbound PTR RR based DNS query request packet traffic through January 1st to December 31st, 2011 (Figure 10).

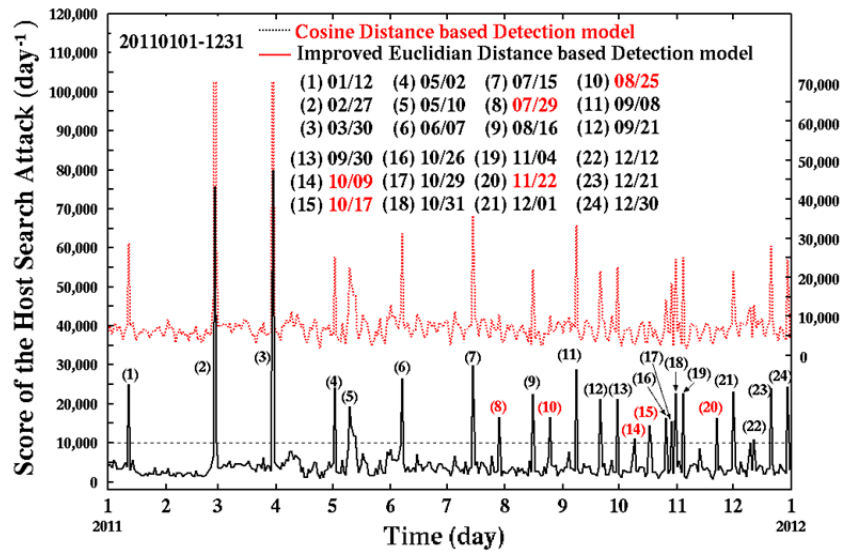In Figure 10, we can observe twenty four peaks (1)-(24) that are assigned  to  (1)  January

Figure 10   Changes in score of the improved Euclidian distance based detection of the source IP address distributed host search (HS) attack in the total PTR resource records (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server through January 1st to December 31st, 2011 (day$^{-1}$ unit).

12th, (2) February 27th, (3) March 30th, (4) May 2nd, (5) 10th, (6) June 7th, (7) July 15th, (8) 29th, (9) August 16th, (10) 25th, (11) September 8th, (12) 21st, (13) 30th, (14) October 9th, (15) 17th, (16) 26th, (17) 29th, (18) 31st, (19) November 4th, (20) 22nd, (21) December 1st, (22) 12th, (23) 21st, and (24) 30th, 2011, respectively.

   Interestingly, in Figure 10, we can observe new five score peaks (8), (10), (14), (15), and (20).   This result shows that the newly suggested detection algorithm has a possibility to increase detection rate.

## 4.  Conclusions

   We developed and evaluated the Euclidean- and cosine-distance based detection models of the source IP address distributed host search (HS) attack in the total inbound PTR resource record (RR) based DNS query request packet traffic through January 1st to December 31st, 2011.   The following interesting results are found: (1) we observed nineteen peaks for host search (HS) attacks in the score changes of the cosine distance based HS attack detection model, however, (2) we found the only twelve peaks in the score changes of the conventional Euclidian based HS attack detection model, (3) we observed the source IP addresses were not

fixed but changed periodically in the newly found peaks, (4) we investigated frequency distribution of the source IP addresses in the newly found peaks, and (5) we developed and evaluated the improved Euclidian distance detection model, resulting that we observed twenty four peaks in the score changes of it.   These results show that the cosine distance based detection model can detect the source IP address distributed HS attack and it has a possibility that the conventional Euclidian distance based detection model also can detect more precisely when taking into consideration the source IP addresses distribution.

   We continue further investigation and development of the HS detection technology in the near future.

## References

1)   Barford, P. and Yegneswaran, V.: An Inside Look at Botnets, Special Workshop on Malware Detection, Advances in Information Security, Springer Verlag, 2006.
2)   Nazario, J.: Defense and Detection Strategies against Internet Worms,I Edition; Computer Security Series, Artech House, 2004.
3)   Kristoff, J.: Botnets, North American Network Operators Group (NANOG32), Reston, Virginia (2004), http://www.nanog.org/mtg-0410/kristoff.html
4)   McCarty, B.: Botnets: Big and Bigger, IEEE Security and Privacy, No. 1, pp.87-90 (2003).
5)   Musashi, Y., Hequet, F., Ludeña Romaña, D. A., Kubota, S., and Sugitani, K.: Detection of Host Search Attacks in PTR Resource Record DNS Query Packet Traffic, IPSJ SIG Technical Reports, Internet Operation and Technology 11th (IOT11) Vol. 2010-IOT-11, No. 8, pp.1-6 (2010).
6)   BIND-9.3.6-P1: http://www.isc.org/products/BIND/
7)   Ludeña Romaña, D. A., Kubota, S., Sugitani K., and Musashi, Y: DNS Based Spam Bots Detection in a University, Proceedings of the First International Conference on Intelligent Networks and Intelligent Systems (ICINIS 2008), Wuhan, China, pp.205-208 (2008).
8)   Lei, M., Musashi, Y., Ludeña Romaña, D. A., Takemori, K., Kubota, S., and Sugitani, K.: Detection of Host Search Activity in Domain Name Reverse Resolution Traffic, IPSJ Symposium Series (IOTS2009), Vol. 2009, No. 15, pp.91-94 (2009).
9)   Musashi, Y., Ludeña Romaña, D. A., Kubota, S., and Sugitani, K.: Detection of Host Name Harvesting Attack in PTR Resource Record Based DNS Query Packet Traffic, IPSJ SIG Technical Reports, Internet Operation and Technology 9th (IOT09) Vol. 2010-IOT-9, No. 9, pp.1-6 (2010).
10)   Musashi, Y., Hequet, F., Ludeña Romaña, D. A., Kubota, S., and Sugitani, K.: Detection of Host Search Activity in PTR Resource Record Based DNS Query Packet Traffic, Proceedings for the Sixth International Conference on Information and Automation (ICIA2010), Harbin, Heilongjiang, China, pp.1284-1288 (2010).