

定点観測システム ISDAS を用いた不正ホスト数の同定

菊池 浩明† 福野 直弥† 寺田 真敏†† 菊地 大輔¶ 土居 範久¶

† 東海大学工学研究科情報理工学専攻
259-1292 平塚市北金目 1117

kikn@tokai.ac.jp, fukuno@ep.u-tokai.ac.jp

†† 日立製作所 Hitachi Incident Response Team (HIRT)

212-8567 神奈川県川崎市幸区鹿島田 890 日立システムプラザ新川崎

¶ 中央大学理工学部情報工学科

112-8551 東京都文京区春日 1-13-27

あらまし 本研究では有限責任中間法人 JPCERT/CC 定点観測システム (ISDAS) のスキャンデータに基づいてインターネット上にワームやウイルスに感染した不正ホストが何台存在するか推定を行う。

Estimation of a Number of Malicious Hosts Based on the Scan Logs Observed by Internet Scan Data Acquisition System

Hiroaki Kikuchi† Naoya Fukuno† Masato Terada†† Daisuke Kikuchi¶
Noriyoshi Doi¶

† Course of Information Engineering, Graduate School of Engineering Tokai University
1117 Kitakaname, Hiratsuka, Kanagawa 259-1292

†† Hitachi, Ltd. Hitachi Incident Response Team (HIRT)

890 Kashimada, Kawasaki, Kanagawa 212-8567

¶ Dept. of Info. and System Engineering, Faculty of Science and Engineering, Chuo University
1-13-27 Kasuga, Bunkyo, Tokyo 112-8551

Abstract In this paper, based on the scan logs observed by JPCERT/CC Internet Scan Data Acquisition System (ISDAS), we try to exactly estimate a number of malicious hosts which are infected with worm or computer virus in the Internet.

1 はじめに

インターネットにはワームやウイルスに感染した不正ホストが多数存在する。不正ホストの代表的な攻撃としてポートスキャンがある。この攻撃は脆弱性のあるホストをランダムに探しては特定のポートへのコネクション確立要求を繰り返す。ワームの代表的なものとして Sasser

がある。[3]によれば、32%の割合で攻撃対象のアドレスを完全にランダムに、23%で1オクテット以外をランダムに生成しては、TCPコネクションの確立を試みる。これら不正ホストの集合について、その大きさや種類、増加割合などの情報はこれまで未知であった。これらの情報はネットワーク管理者にとって有効であり、

効率的な安全対策を施すことに応用できる。

これまでに我々は3台のセンサで3ヶ月間観測した結果に基づいての不正ホストの総数を同定した[1]。しかし、センサ数や観測期間に関してデータ不足があり、その推定結果の信頼性に問題があった。そこで本稿では、有限責任中間法人 JPCERT コーディネーションセンター (以後、JPCERT/CC) の定点観測システム ISDAS のデータについて、[1]と同じ手法を適用して推定した不正ホストの母集団の性質について報告する。

2 潜在不正ホスト数の同定問題

2.1 モデル

不正ホストとは、ウィルスやネットワークワームなどにより他のホストへの攻撃(ポートスキャン)を仕掛けるホストである。センサとは、不正ホストからの攻撃を観測する正規ホストであり、決して感染しない。有効な全グローバルアドレスの数を n_0 、不正ホストの数を n とする。 x 台の独立したセンサで期間 $[0, t]$ で観測できる異なるセンサの累積数をユニークホスト数と呼び、 $h(x, t)$ で表す。

潜在不正ホスト数の同定とは、複数の分散センサで観測されたユニークホスト数 $h(x, t)$ を与えて、全不正ホスト数 n を同定する問題である。センサは複数箇所に分散して設置しているが、不正ホストはその場所を知らない。

この問題を、次を仮定する単純化されたモデルで考えよう。

仮定1 静的アドレス不正ホストは静的なアドレスを持つ。(DHCPなどによる動的なアドレス割り当ては行なわない)。IPアドレスの詐称は行なわない。

仮定2 センサの独立性 スキャン先はランダムに決定し、有効なアドレス空間を一様分布する。ポートの区別は行なわない。

仮定3 観測時期・期間の独立性 不正ホストの振る舞いは常に一定であり、観測時期や期間の長さには依存しない。

仮定4 不正ホストの独立性 全ての不正ホストは単位時間当たり c 回のスキャンを実行する。スキャンの割合は時刻にも不正ホストにも依らず一定とする。

このモデルの上では、あるセンサが攻撃対象に選ばれる確率は $1/n_0$ であり、実際には n 台の不正ホストがあるので単位時間にスキャンを受ける確率は n/n_0 といえる。仮定4により、単位時間にセンサが観測する平均ユニークホスト数 a は、

$$a = c \frac{n}{n_0} \quad (1)$$

で与えられる。

センサ数 x や観測期間 t を増やしても、ユニークホスト数はそれらに対して線形ではなく、やや鈍い増加を示すはずである。このゆがみは、不正ホストの総数 n に依存して大きくなり、それゆえ、このゆがみを正確に測定できればそこから n が求められるだろう。

ここで、 n の上限は 2^{32} の全IPアドレス空間であるが、実際には未割り当てや外部へ公開していないプライベートなアドレスブロックがある¹。

2.2 観測期間についてのユニークホスト数

仮定より、 $h(1, 1) = a$ であり、次の単位時刻には更に a 台の不正ホストが観測できる。センサ数 $x = 1$ とおいて、 $h(t) = h(1, t)$ と置き換え、一般化すると、

$$h(t+1) = h(t)(1 - a/n) + a$$

が得られる。差分 $h(t+1) - h(t)$ から極限を取り、ユニークホスト数に対する微分方程式

$$\frac{dh}{dt} = -\frac{a}{n}h(t) + a \quad (2)$$

が得られる。この微分方程式を解くと

$$h(t) = n(1 - e^{-\frac{a}{n}t}) \quad (3)$$

を得る。ここで、 n が潜在的な不正ホストの数、 a が単位時間にセンサが観測する平均ユニークホスト数である。

¹[6]によると、2005年7月時点で $n_0 = 353,284,184$ が示されている。

2.3 センサ台数についてのユニークホスト数

前節では、 $h(1, t) = h(t)$ と置いたが、 t を x と置き換えても全く同様な議論が成立する。すなわち、センサ数 x についてのユニークホスト数もまた式 (4) で定式化できる。

$$h(x) = n(1 - e^{-\frac{a}{n}x}) \quad (4)$$

従って、観測期間を変えてのユニークホスト数から導く不正ホスト数と観測センサ数についてのユニークホスト数から推定する不正ホスト数とが一致することが期待できる。

3 解析

3.1 定点観測データ ISDAS

ISDAS(Internet Scan Data Acquisition System) は、JPCERT/CC が運用しているセキュリティに関するトラフィックの分散観測システムである [5]。単位時間当たりの主要なポート (13, 445, 135, 139, 1026, 80) のトラフィックなどを毎日更新している。

2 節で説明した不正ホスト数を同定するためには十分な観測期間における独立した複数のセンサの観測データを必要とする。そこで、JPCERT/CC に 2004 年 9 月 1 日～2005 年 9 月 30 日の間の独立した 12 台のセンサのログデータを提供してもらった。

3.2 解析方法

観測したスキャンデータに 2 節のモデルを最小二乗法であてはめを行い、パラメータ n と a を同定する。観測条件を明確にするために、ユニークホスト数を次の $h(S, T, P)$ で示す。

S: センサ集合 ($S = \{s_1, \dots, s_{13}\}$)

T: 観測期間

P: 観測ポート番号

例えば、 $h(9, [2004/5 - 2004/7], 135)$ は s_9 におけるポート 135 の 2004 年 5 月から 3 ヶ月間の

累積のユニークホスト数である。同様に、ここから算出された潜在不正ホスト数を $n(S, T, P)$ と表す。本報告では次の解析を行った。

1. 観測期間からの潜在不正ホスト数同定。式 (3) に基づき、12 台のセンサ各々について、代表的なポート 135, 139, 445 毎の総不正ホスト数を同定する。
2. センサ数の累積からの潜在不正ホスト数の同定。式 (4) に基づき、センサ数 x を 1 から 12 まで変化させて、不正ホスト数を同定する。ただし、センサによってスキャン数は一定でないので、累積する順序によってユニークホストの増加量が変わる。そこで、スキャン数についてセンサをソートし、昇順と降順の平均を取る。
3. 観測期間の独立性の検証。仮定 3 の妥当性をみるために、観測データを 3 ヶ月に分割し、 n の変動を調べる。
4. センサ独立性の検証。仮定 2 を確かめるため、ユニークホストだけではなく、センサ毎のスキャン数を調べる。スキャン数とユニークホスト数の関係を明らかにする。加えて、全てのセンサの組み合わせについての累積ユニークホスト数を求め、センサ間の相関を調査する。
5. 不正ホストの独立性の検証。全センサで観測される渡り歩きの活発な不正ホストや特定の不正ホストに集中している執着不正ホストがどのくらいいるのか、各不正ホストのソースアドレスについて、観測されたセンサの数を調べる。

3.3 観測期間からの不正不正ホスト数同定

T_1 を 2004 年 9 月 1 日から 2005 年 9 月 30 日までの約 1 年間とする。表 2 にセンサ s のスキャン数、ユニークホスト数 $h(1, T_1, 445)$ 、一日当りのユニークホスト数の増加量 $\Delta h(s, T_1, 445)$ を示す。最小のユニークホストをもつ s_9 と最大の s_1 の間に約 10 倍の差がある。図 1 に一日あたりのユニークホスト数の増加数 $\Delta h(s, T_1, 445)$ を

表 2: 各センサの観測データ

センサ	総スキャン数	$h(x)$	$\Delta h(x)/[\text{日}]$
s_1	268024	97102	245.8
s_2	153310	63198	160.0
s_3	154126	60755	153.8
s_4	137848	40315	102.1
s_5	168191	62881	159.2
s_6	173566	47809	121.0
s_7	17167	10066	25.5
s_8	164078	54865	138.9
s_9	10667	9046	22.9
s_{10}	170417	24394	61.8
s_{11}	30898	13200	33.4
s_{12}	143725	53716	136.0

示す。観測期間を広げていくことによって、増加量がゆるやかに減っている。

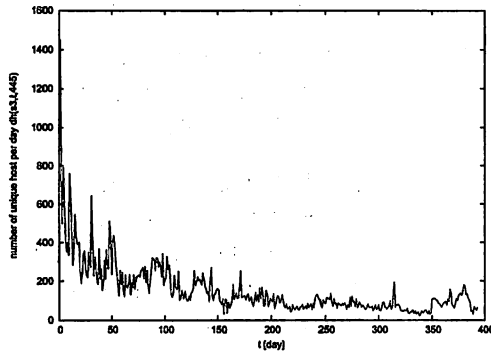


図 1: ユニークホスト数の一日当たりの増加数 $\Delta h(s_3, \Delta t, 445)$

これらのデータについて、あてはめを行った結果を図 2 に示す。代表的なセンサ s_1, s_3, s_{11} について図示している。いずれも、式(1)の理論式と実測値の間に一致が見られる。全てのセンサについての推定不正ホスト数 $n(s, T_1, 445)$ を単位時間当たりのスキャン実行数 a を表 1 に示す。

センサ s_7 と s_9 については他と比べて誤差が 3 桁ほど大きい。この 2 センサのユニークホストを図 3 に示す。図より、観測期間の途中から急激にユニークホストが増加し、 t に対して $h(x)$ が線形に増加していることがわかる。そこでこれらを除外して平均を求めている。同定された不正ホスト数の確率分布を図 5 に示す。

図 4 に、ポート $P=135, 139, \text{ICMP}$ の各々に

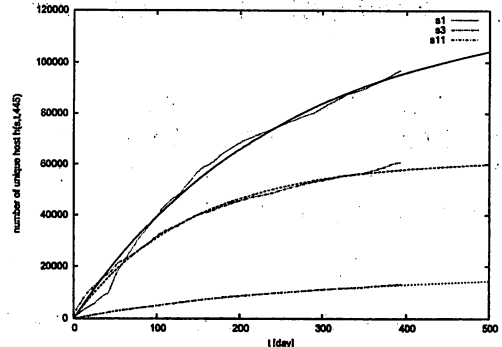


図 2: ユニークホスト数の推移 $h(s, t, 445)$

についての累積ユニークホスト数 $h(s_3, T_1, P)$ を表す。図より、総量は各々異なるが、ユニークホスト数の増加については同様の振る舞いをしていることが観測できる。そこで、以降の議論は代表的なポート 445 についてのみ行うこととする。他のポートについても、同様の結果が導かれると考えられる。

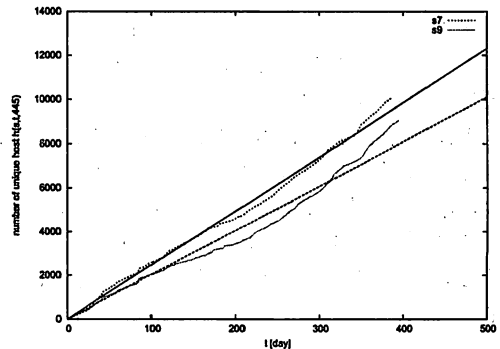


図 3: 当てはめの失敗例

3.4 センサ台数 x からの不正ホスト数同定

図 6 は観測期間 $T_2 = [2005/2/1 - 2005/2/28]$ におけるセンサ台数についての累積ユニークホスト数 $h(x, T_2, 445)$ を示している。12 台のセンサを、表 2 の総スキャン数について並べ替え、昇順と降順に対応する。最大値と最小値を求め、その平均値 $\bar{h}(x, T_2, 445)$ に対して、式(4)のあてはめを実行した。このとき同定された潜在不正ホスト数 $n(12, T_2, 445)$ を表 3 に示す。

表 1: 1 年間の観測期間から推定した総不正ホスト数 $n(1, T_1, 445)$

s	n	誤差 [%]	n/a	誤差 [%]	c [回/秒]
s_1	121375	1.02	255.769	1.77	15.99
s_2	76925.4	0.82	245.321	1.46	16.67
s_3	61939.2	0.34	145.627	0.83	28.08
s_4	49056.1	0.60	249.579	1.03	16.38
s_5	68167.2	0.25	170.927	0.53	23.92
s_6	58291	0.89	242.368	1.59	16.87
s_7	4.59E+09	8.70E+03	1.87E+08	8.71E+03	0.00
s_8	64973.4	0.61	239.498	1.10	17.07
s_9	1.11E+09	6.56E+03	5.51E+07	6.57E+03	0.00
s_{10}	30669.7	0.80	262.941	1.37	15.55
s_{11}	17562.5	0.45	297.874	0.72	13.73
s_{12}	75299.2	1.09	330.488	1.69	12.37
平均 \bar{n}_1	62425.8	0.68	244.0392	1.21	16.75

表 3: 観測期間から推定した総不正ホスト数 $n(S, T_1, 445)$

観測時期 T_2	n_2	誤差 [%]	n_2/a	誤差 [%]	c [回/秒]
2005/2/1 - 2005/2/28	111655	24.96	33.113	28.76	123.48

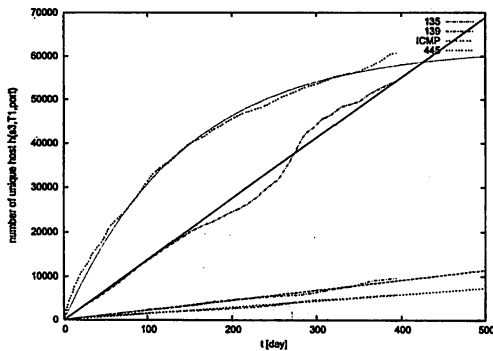


図 4: ポートによるユニークホスト数 $h(s_3, T_1, port)$

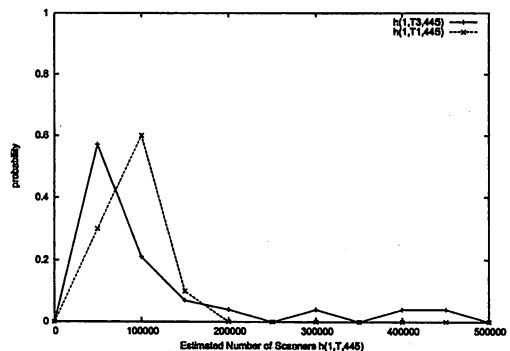


図 5: 推定不正ホストの確率分布

3.5 観測期間についての独立性

表 4 は、観測期間による変動をみるために、2004 年 9 月から 2005 年 7 月までを 3 ヶ月ごとに分割した区間 T_3 についての潜在不正ホスト数同定結果を示している。 n が 10^6 を超えたものは、あてはめに失敗したとみなした。(特に、センサ s_7 については、全区間について失敗している。) センサと観測期間の組は $12 \times 4 = 48$ 組あるが、あてはめが有効だったものは 28 組だけであり、あてはめの成功率は 50% であった。これらの不正ホスト数 n_3 の確率分布を図 5 に

示す。平均 $\bar{n}_3 = 87735$ であり、1 年間の観測期間 T_1 から求めた $\bar{n}_1 = 62425$ に近い値が同定されている。

十分に短い区間で見れば、スキャンはランダムに実行されている。図 7 はセンサ s_3 におけるスキャンの到着間隔の分布を表している。理想的なポアソン分布に近い分布をしていることが裏づけされる。

3.6 センサの独立性

表 2 にセンサごとのスキャン数とユニークホスト数を示す。スキャン数では最小の s_9 と最大

表 4: 3ヶ月の観測期間から推定した総不正ホスト数 $n(1, T_2, 445)$

観測開始時期	2004/09	2004/12	2005/03	2005/06	平均 \bar{n}_3
s_1	3.76E+09	95341.5	43903.7	111826	83690.40
s_2	2.44E+08	20692.4	376112	9.56E+06	198402.20
s_3	44780.3	32577.9	38465.2	32531.8	37088.80
s_4	199272	5.46E+08	25017.3	1.24E+08	112144.65
s_5	92525.3	32994.4	25017.3	1.24E+08	50179.00
s_6	55772.7	83206.4	1.30E+09	72061.1	70346.73
s_7	7.50E+08	3.42E+07	3.23E+10	2.71E+07	0
s_8	426469	22085.6	136266	1.47E+08	194940.20
s_9	2.04E+06	9750.29	1.69E+08	3.89E+07	9750.29
s_{10}	13579	1.57E+08	1.13E+08	1.70E+08	13579.00
s_{11}	31652.7	28459.9	39735.4	7953.8	26950.45
s_{12}	2.37E+08	96056.9	1.40E+10	262474	179265.45

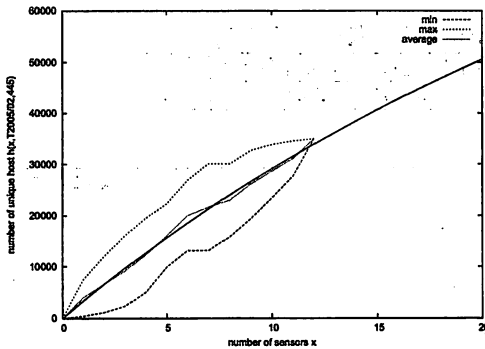


図 6: センサ台数についての累積ユニークホスト $h(x, T_2, 445)$

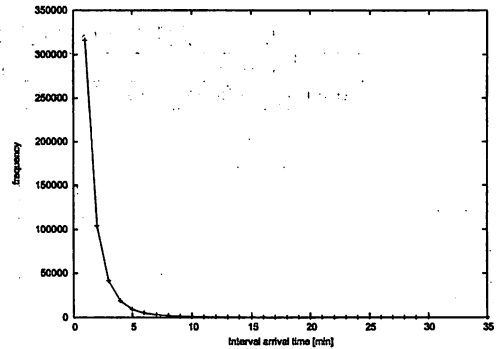


図 7: 到着間隔

の s_1 の間には約 20 倍の差がある。

図 8 はポート 445 についての観測アクセス回数と $h(x, T_1, 445)$ の関係を示す。両者の相関係数は 0.93 であり、明らかな正の相関が示されている。

センサが十分に独立していることを確かめるために図 9 に全ての異なる 2 組のセンサによる累積ユニークホスト数を示す。図の等高線は

$$r_{(i,j)} = \frac{h(s_i, T_1) + h(s_j, T_1) - h(\{s_i, s_j\}, T_1)}{h(\{s_i, s_j\}, T_1)} \quad (5)$$

で定義される重複度 $r_{i,j}$ 値を示している。ただしここで、 $h(S, T)$ はポート 445 についてのセンサ集合 S の累積ユニークホスト数とする。図より、 s_4 と s_8 の間、及び、 s_4 と s_6 間、 s_7 と s_9 間の 3 箇所の相関が強いことが示されている。

3.7 不正ホストの独立性

図 10 に異なるスキャン先センサ数についての不正ホスト (ソースアドレス) 数の分布を示す。図より、最も多いのは、単一のセンサにのみ観測された不正ホストであり約 100 万個 (86%) 存在し、12 台全てにスキャンを行ったものは 293 個 (0.02%) であった。一台の不正ホストは年間平均 1.2 台のセンサにスキャンしている。ただし、ここでは、全ポートの値を総和して

図 11 は、不正ホストのソースアドレスの IP アドレス空間に置ける分布を表している。アドレスの第一オクテットのみに着目し、1 つでもスキャンが観測できた 24 ビットのアドレスブロックをアクティブとみなして図示している。明らかに、すべてのアドレスブロックから一様にスキャンが届いていることが示されている。

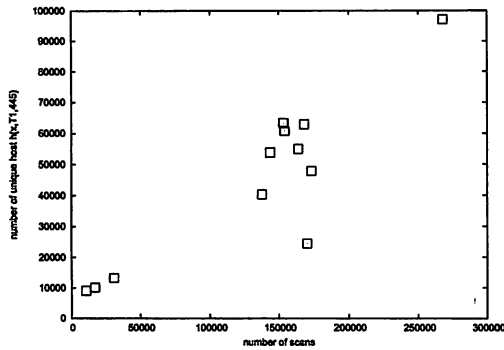


図 8: 総アクセス数とユニークホスト数

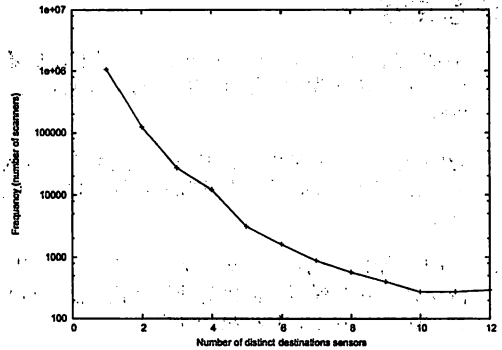


図 10: センサ先の数についての不正ホスト数

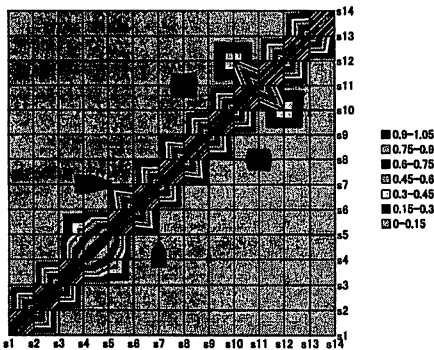


図 9: 観測センサの重複度合い $r_{i,j}$

4 結論

ポート 445 を撃つ不正ホストの数は、1 年間の観測期間から推定すると $n_1 = 8773 (\pm 36000)$ 、3 ヶ月の複数の区間からは $n_3 = 87735 (\pm 21000)$ 、12 台の独立したセンサからは $n_2 = 111655$ であった。ただし、誤差は 95% (2σ) の信頼区間による。従って、8 万台を中心に分布しており、高々 30 万台であることが示された。また、これらの不正ホストは T_1 を平均して毎秒 16.75 回のスキャンを実行していた。

モデルの妥当性を確かめるために、仮定 2, 3, 4 に相当する 3 種の独立性を検証した。当てはめに成功した区間については、観測期間に依存せず、一定のホスト数が推定できたが、成功確率は 50% であった。

12 台のセンサにはスキャン数で約 20 倍の差

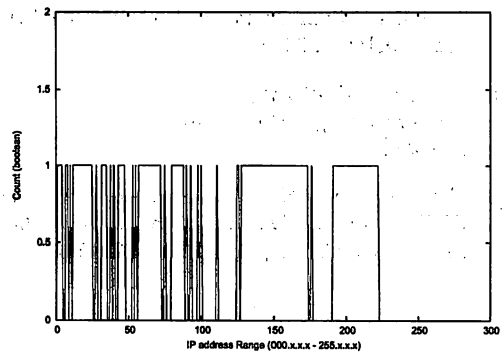


図 11: スキャン元の第一オクテットの分布

があるが、ユニークホスト数で見ると、全 66 組中 3 組を除いて大きな相関はなく、十分独立していた。スキャン数とユニークホストには正の相関があった。不正ホストは全アドレス空間上に一様に分布しており、活発の差はあるが年間平均 12 台中の 1.2 台のセンサ (図 10) で観測されていることがわかった。

今後の課題としては、モデルの仮定の一般化、あてはめの失敗の適切な処理が挙げられる。

謝辞

本研究を遂行するにあたり、定点観測データ提供して頂いた JPCERT/CC に感謝する。本研究に関して有意義なご助言を頂いた竹田 春樹氏、中谷 昌幸氏、鎌田 敬介氏 (JPCERT/CC) に感謝する。

参考文献

- [1] 菊池 他, ネットには何台の不正ホストがいるのか?, 情報処理学会, CSS 2005, pp.421-426, 2005.
- [2] 杉山 他, アクセスログを用いた不正ホスト総数の推定に関する検討, 情報処理学会, FIT 2005, 2005.
- [3] 寺田, 高田, 土居, ネットワークワーム動作検証システムの提案, 情報処理学会論文誌, Vol. 46, No. 8, pp. 2014-2024, 2005.
- [4] J. Jung, V. Paxson, A. W. Berger, and H. Balakrishnan, "Fast Portscan Detection Using Sequential Hypothesis Testing", proc. of the 2004 IEEE Symposium on Security and Privacy (S&P'04), 2004.
- [5] JPCERT/CC,ISDAS,
(<http://www.jpCERT.or.jp/isdas>, 2006 年
2月参照)
- [6] "Number of Hosts advertised in the DNS", Internet Domain Survey, July,2005.
(<http://www.isc.org/ops/reports/2005-07>).