

洗練されたサイバー攻撃に対する 最適な欺瞞的防御戦略の選択法

大阪府立大学大学院 工学研究科

浜崎 晃次 (Koji Hamasaki), 北條 仁志 (Hitoshi Hohjo)

Graduate School of Engineering, Osaka Prefecture University

1 はじめに

昨今、情報技術の急速な発展に伴い、ネットワークは人間社会に大きな利便性をもたらしており、様々な分野で非常に幅広い応用がなされている。しかし、ネットワークセキュリティの問題はますます深刻化しており、様々なサイバー攻撃がネットワークに深刻な脅威を与えている。サーバーやネットワークなどのリソースに過剰な負荷をかけたり、脆弱性を突いたりすることでサービスを妨害する DoS 攻撃や、特定の組織や個人に狙いを定め、それに適した攻撃を組み合わせる APTs(持続的標的型攻撃) 等がその例として挙げることができる。

これらの攻撃の中でも、我々は DoS 攻撃に焦点を当てる。DoS 攻撃は前述したようにサービスを妨害して、正当なユーザーがリソースにアクセスできないようにする攻撃である。攻撃者はネットワーク全体を標的にして、一時的または恒久的な利用不能を引き起こす。分散型 DoS (DDoS) 攻撃では、防ぐことも回復することもより困難になる。また攻撃者は、正当なユーザーになりすまして防御者の情報を得ることができるため、防御者よりも時間的にもステルス性にも優れている。そこで、この防御者にとって不利な状況を解決するために、欺瞞的防御を利用することができる。欺瞞的防御は、攻撃者が組織内部ネットワークに侵入することを前提としており、ハニーポット等の囮となるツールをネットワーク内部に設置することで攻撃者を混乱させ、攻撃を遅らせることを可能とする。その間に、防御者は攻撃に対する解決策を立てることができる。この欺瞞的防御を利用するために防御者は追加のコストを費やすことになるが、この技術は DoS 攻撃を軽減する有効的な手段である。

このような欺瞞的防御の効果を分析するためにゲーム理論を用いて様々な研究がなされてきた。2011 年、Carroll et al.[2] はシグナリングゲームを用いて、ハニーポットを利用した防御者と攻撃者の間の相互作用をモデル化した。通常のシステムをハニーポットとして偽装したり、ハニーポットを通常のシステムとして偽装したりする欺瞞的防御が、防御者の最適戦略になり得ることを示した。2016 年には、この研究を Ceker et al.[3] が DoS 攻撃を含むより広い局面に拡張し、利益と損害を定量化するメカニズムを提案した。しかし、これらの研究においては防御者と攻撃者の間の相互作用を 1 段階のシグナリングゲームで捉えた簡単なモデルであった。

本論文では、ハニーポットを利用した防御者と攻撃者の間の相互作用を多段階シグナリングゲームで捉え、より現実の状況へと近づけたモデルを検討する。さらに、防御者の欺瞞信号の影響を定量化するために、欺瞞信号の減衰係数を考慮する。これによって、長期的な DoS 攻撃と防御の動的分析と推論を実現する。本研究

成果は、DoS 攻撃と防御の対立を研究するための効果的なモデル化手法を提供し、ネットワークセキュリティの分野で欺瞞的防御を適用するための理論的な指針となると考えられる。

2 シグナリングゲーム

シグナリングゲーム [4] は Dos 攻撃者と防御者の 2 人のプレイヤーによって行われる。防御者は不正なアクセスを検知するための罫であるハニーポットを配置してネットワークを守ることを目的としている。また、費用対効果を高めるためにシステムを偽装することができる。ネットワークを構築した後、攻撃者はシステムを侵害しようとする。攻撃者は、通常のシステムを侵害することは可能であるが、ハニーポットを侵害することは不可能である。攻撃者がハニーポットを侵害しようとした場合、防御者はその行為を観察し、後に防御を改善することができる。このような防御者と攻撃者の相互作用をシグナリングゲームとしてモデル化する。

2.1 仮定

DoS（特に DDoS）攻撃は通常、大量のコンピュータによって行われるが、ここでは単一の集中型攻撃者がサーバーに対して過剰なアクセスやデータを送信し、利用不能を引き起こす場合に限定している。したがって、分散型攻撃者（複数のハッカーグループ、国、企業など）の場合は対象外とする。

ゲーム中、攻撃者は防御者が送信した信号を観測した後、防御者のタイプに関する知識を更新することができる。しかし、ここでは簡単のために、他の種類の観測（スパイ攻撃やプローブ攻撃など）は含めないことにする。最後に、プレイヤーは完全に合理的であり、自分の効用を最大化したいと考えていると仮定する。

2.2 シグナリングゲームモデルの定義

本研究におけるシグナリングゲームモデルのルールを以下に記す。

1. 自然が防御者のタイプを通常のシステムまたはハニーポットのいずれかに決定する。
2. 防御者は自身のタイプに基づいて、攻撃者に対して信号を送信する。
3. 攻撃者は信号を受け取り、攻撃、観察、退却の中から行動を選択する。
4. 防御者のタイプと信号、攻撃者の行動に基づいて、両プレイヤーに利得を与える。

上記の 1~4 を行い、ゲームは終了する。

防御者から攻撃者に送信される信号は、「実信号」と「欺瞞信号」の 2 種類に分けられる。

定義 1 実信号とは、必然的に公開された防御者のプライベートな情報のことを表す。この信号を送るために、コストが発生することはない。

定義 2 欺瞞信号とは、本当のタイプを隠すために送る信号のことで、攻撃者の判断を誤らせるように誘導することができる。理由なく信号が発生することはないため、欺瞞信号を送信するためには、追加のコストを支

払う必要がある。

これにより、本研究におけるシグナリングゲーム (SG) モデルを以下のように定義する。

定義 3 $SG = (N, T, M, S, P, U)$ とする。

1. $N = (N_D, N_A)$ はゲームのプレイヤー集合である。 N_D は防御者、 N_A は攻撃者である。
2. $T = (t_D, t_A)$ はタイプ集合である。 T_D は防御者のタイプであり、 $T_D = (t_N, t_H)$ とする。 t_N は通常のシステム、 t_H はハニーポットである。 T_A は攻撃者のタイプであり、 $T_A = (DoS \text{ 攻撃者})$ とする。
3. $M = (m_N, m_H)$ は防御者の信号集合である。 m_N は自身を通常のシステムであると思わせる信号、 m_H は自身をハニーポットであると思わせる信号である。
4. $S = (A, O, R)$ は攻撃者の行動集合である。 A は攻撃、 O は観察、 R は退却である。
5. $P = (P_A, \tilde{P}_A)$ は攻撃者の信念集合である。 P_A は事前確率、 \tilde{P}_A は事後確率である。
6. $U = (U_D, U_A)$ は防御者と攻撃者の利得関数集合である。 U_D は防御者の利得関数、 U_A は攻撃者の利得関数である。

図 1 は、シグナリングゲームのイメージ図である。また、本研究モデルで使用する表記法を表 1 にまとめた。

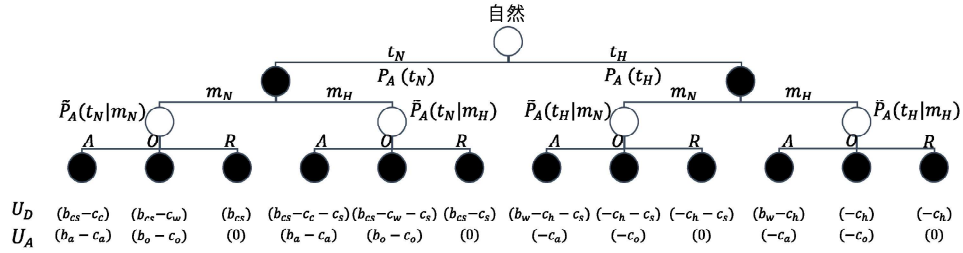


図 1 シグナリングゲーム

2.3 利得の定量化

本節では、攻撃者が防御者の資産評価とコストについて不確実であるシナリオを考える。まず、DoS 攻撃のコストを防御者側と攻撃者側に分けて定量化する。

Basagiannis et al.[1] は、Meadows[6] のフレームワークを用いて、DoS 攻撃の防御者側と攻撃者側のコストを定量化する確率モデルを提案した。このモデルを参照し、本研究モデルでは $c_c = 4000$ 、 $c_a = 600$ 、 $c_w = 80$ 、 $c_o = 30$ とする。これらの値を当てはめて、防御側のサービスレートを推定し、攻撃者の戦略による劣化を分析できるようにする。有効なサービスレートに対する顧客の満足度を測定するために、以下の式を使用する。

$$U(R) = 0.16 + 0.8 \cdot \ln(R - 3) \quad (1)$$

式 (1) は、防御者がサービスレート R でサービスを提供したりソースにアクセスした顧客の満足度を定量化したものである。攻撃者の目的は、顧客がリソースを利用できない状況にすることなので、攻撃によってレート

表 1 表記法

表記	説明
Defender	
T_N	通常のサーバー
T_H	ハニーポット
m_N	自身を通常のサーバーであると思わせる信号
m_H	自身をハニーポットであると思わせる信号
c_c	通常のサーバーが攻撃されることによるコスト ($c_c \geq 0$)
c_s	欺瞞信号を送信することによるコスト ($c_s \geq 0$)
c_h	ハニーポットを設置することによるコスト ($c_h \geq 0$)
c_w	通常のサーバーが観察されることによるコスト ($c_w \geq 0$)
b_{cs}	通常のサーバーに対する顧客の満足度
b_w	ハニーポットで攻撃者を観察することによる利益 ($b_w \geq 0$)
R_d	サービスレート
Attacker	
A	攻撃
O	観察
R	退却
$P_A(t_N)$	標的が通常のサーバーである確率
$P_A(t_H)$	標的がハニーポットである確率
$P_A(t_N m_N)$	m_N が送られてきた時に標的が通常のサーバーである確率
$P_A(t_N m_H)$	m_H が送られてきた時に標的が通常のサーバーである確率
$P_A(t_H m_N)$	m_N が送られてきた時に標的がハニーポットである確率
$P_A(t_H m_H)$	m_H が送られてきた時に標的がハニーポットである確率
c_a	攻撃コスト ($c_a \geq 0$)
c_o	観察コスト ($c_o \in [0, c_a]$)
b_a	攻撃による利益 ($b_a \geq c_a$)
b_o	観察による利益 ($b_o \geq c_o$)
R_a	攻撃レート
R_o	観察レート

1. 最適な攻撃戦略 $S^*(m)$ を選択する。

$$S^*(m) \in \max_{m \in M} \sum P_A(T_D|m) U_A(m(T_D), S, T_D) \quad (5)$$

2. 最適な防御戦略 $m^*(T_D)$ を選択する。

$$m^*(T_D) \in \max_{m \in M} \sum U_D(m, S^*(m), T_D) \quad (6)$$

3. 事後確率をベイズルールによって改良する。

$$\tilde{P}_A(T_D) = \tilde{P}_A(T_D|m) \quad (7)$$

3 多段階シグナリングゲーム

通常のシグナリングゲームは、防御者と攻撃者の1回の対立のみを捉えている。現実世界の DoS 攻撃では、防御者は長期的に同様の独立した相互作用に複数回直面する可能性があると考えられる。このような状況を、シグナリングゲームを繰り返し行う多段階シグナリングゲームでモデル化することができる。また、Zhang et al.[7] は防御者の防御信号の影響を定量化するために、信号減衰係数を用いている。本研究でも信号減衰係数を用い、長期的な DoS 攻撃と防御の対立の中で欺瞞信号の効果が減衰していく状況を捉える。

3.1 多段階シグナリングゲームモデルの定義

本研究における多段階シグナリングゲーム (*MSSG*) モデルを以下のように定義する。

定義 4 $MSSG = (N, K, T, M, S, \delta, P, U)$ とする。

1. $N = (N_D, N_A)$ はゲームのプレイヤー集合である。 N_D は防御者、 N_A は攻撃者である。
2. K はゲームステージの総数である。 $G(k)$ は k ステージの *SG* を表す。また、 $k = \{1, 2, \dots, K\}$ 。
3. $T = (t_D, t_A)$ はタイプ集合である。 T_D は防御者のタイプであり、 $T_D = (t_N, t_H)$ とする。 t_N は通常のシステム、 t_H はハニーポットである。 T_A は攻撃者のタイプであり、 $T_A = (DoS \text{ 攻撃者})$ とする。
4. $M = (m_N, m_H)$ は防御者の信号集合である。 m_N は自身を通常のシステムであると思わせる信号、 m_H は自身をハニーポットであると思わせる信号である。
5. $S = (A, O, R)$ は攻撃者の行動集合である。 A は攻撃、 O は観察、 R は退却である。
6. δ は欺瞞信号の減衰係数である。あるステージでの欺瞞信号が前回のステージからどの程度減衰しているのかを示す。また、 $0 \leq \delta \leq 1$ 。
7. $P = (P_A, \tilde{P}_A)$ は攻撃者の信念集合である。 P_A は事前確率、 \tilde{P}_A は事後確率である。
8. $U = (U_D, U_A)$ は防御者と攻撃者の利得関数集合である。 U_D は防御者の利得関数、 U_A は攻撃者の利得関数である。

3.2 多段階シグナリングゲームの均衡解

多段階ゲームでは、プレイヤーの最適戦略が各ステージで別々に求められ、全ての最適戦略が多段階ゲームの均衡解を構成する。

第1ステージの均衡解

攻撃者は以前の対決から実際の防御者のタイプを分析することができないため、防御者が送信する欺瞞信号には減衰効果がない。

第2ステージ以降の均衡解

前回のステージで修正した事後確率を事前確率として使用し、自然の役割が置き換えられる。また、ステージを経たことにより攻撃者の信号識別能力が向上し、防御者の欺瞞信号が減衰する。防御者が欺瞞信号を送信した場合は $\delta = 0.7$ とする。実信号を送信した場合は $\delta = 1$ とし、攻撃者は信号識別能力を向上させることはできない。

第 K ステージの均衡解

ゲームが進み、防御者が欺瞞信号を発信する回数が増えると、 $\delta^{k-r-1} = 0$ となる。ここで r は防御者が実信号を送信したステージの数である。この段階で防御者の欺瞞信号が効果をなくし、本当のタイプを偽装することができなくなる。ゲームは不完全情報静的ゲームに変わる。

以上より、全てのステージでそれぞれ均衡解が求められ、全てのステージにおける最適戦略は全体のゲームプロセスにおける最適戦略を構成することになる。

4 解析

多段階シグナリングゲームにおける第1ステージの均衡解を解析する。一段階のシグナリングゲームの均衡は、大きく分けて2つ存在する。一つ目は分離均衡である。この均衡は、防御者のタイプと同じ数だけの信号が存在し、どのタイプもそれぞれ一つの信号が割り当てられている時に存在する。表2に全ての分離均衡をまとめる。

表2 分離均衡

名前	$(m(t_N), m(t_H))$ $-(s(m_N), s(m_H))$	Conditions	μ, γ
E1	$(m_N, m_H) - (A, R)$	$R_d \geq \frac{c_c - c_s}{v_1 - v_2}, R_d \geq \frac{b_w - c_s}{v_3 - v_4}, R_a \leq \frac{c_c}{v_a}$	1, 0
E2	$(m_N, m_H) - (R, R)$	$R_d \leq \frac{c_s}{v_2 - v_1}, R_d \leq \frac{c_s}{v_4 - v_3}, R_a > \frac{c_c}{v_a}$	1, 0
E3	$(m_H, m_N) - (A, R)$	$R_d > \frac{c_c + c_s}{v_2 - v_1}, R_d > \frac{b_w + c_s}{v_4 - v_3}, R_a \leq \frac{c_c}{v_a}$	0, 1
E4	$(m_H, m_N) - (R, R)$	$R_d > \frac{c_s}{v_2 - v_1}, R_d > \frac{c_s}{v_4 - v_3}, R_a > \frac{c_c}{v_a}$	0, 1

左の列から均衡の名前、防御者と攻撃者の戦略の組み合わせ、均衡条件、事後確率である。表にまとめやす

くするために、 $\mu = \tilde{P}_A(t_N|m_N), \gamma = \tilde{P}_A(t_N|m_H)$ としている。E1, E2 では t_N が m_N 、 t_H が m_H を送信し、E3, E4 では t_N が m_H 、 t_H が m_N を送信するという風に、それぞれのタイプが別の信号を送信している状況での均衡である。

二つ目は一括均衡である。この均衡は、防御者が自分のタイプによらず同じ信号を送信する時に存在する。表 3 に全ての一括均衡をまとめる。

表 3 一括均衡

名前	$(m(t_N), m(t_H))$ $-(s(m_N), s(m_H))$	Conditions	Prior&Posterior
E5	$(m_N, m_N) - (A, A)$	$\frac{c_s}{v_2-v_1} \geq R_d \geq \frac{c_s}{v_4-v_3}$	$P_A(t_N) \geq \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, P_A(t_N) \geq \frac{R_o \cdot v_a}{c_c},$ $\gamma \geq \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, \gamma \geq \frac{v_a \cdot R_a}{c_c}$
E6	$(m_H, m_H) - (A, A)$	$\frac{c_s}{v_4-v_3} > R_d > \frac{c_s}{v_2-v_1}$	
E7	$(m_N, m_N) - (A, O)$	$\frac{c_s+c_w-c_c}{v_2-v_1} \geq R_d \geq \frac{c_s-b_w}{v_4-v_3}$	$P_A(t_N) \geq \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, P_A(t_N) \geq \frac{R_o \cdot v_a}{c_c},$ $\gamma < \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, \gamma \geq \frac{v_a \cdot R_o}{c_w}$
E8	$(m_H, m_H) - (A, O)$	$\frac{c_s-b_w}{v_4-v_3} > R_d > \frac{c_s+c_w-cc}{v_2-v_1}$	
E9	$(m_N, m_N) - (A, R)$	$\frac{b_w-c_s}{v_3-v_4} \geq R_d \geq \frac{c_c-c_s}{v_1-v_2}$	$P_A(t_N) \geq \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, P_A(t_N) \geq \frac{R_o \cdot v_a}{c_c},$ $\gamma < \frac{v_a \cdot R_a}{c_c}, \gamma < \frac{v_a \cdot R_o}{c_w}$
E10	$(m_H, m_H) - (A, R)$	$\frac{c_c-c_s}{v_1-v_2} > R_d > \frac{b_w-c_s}{v_3-v_4}$	
E11	$(m_N, m_N) - (O, A)$	$\frac{c_s+c_c-c_w}{v_2-v_1} \geq R_d \geq \frac{c_s+b_w}{v_4-v_3}$	$P_A(t_N) < \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, P_A(t_N) \geq \frac{R_o \cdot v_a}{c_w},$ $\gamma \geq \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, \gamma \geq \frac{v_a \cdot R_a}{c_c}$
E12	$(m_H, m_H) - (O, A)$	$\frac{c_s+b_w}{v_4-v_3} > R_d > \frac{c_s+c_c-c_w}{v_2-v_1}$	
E13	$(m_N, m_N) - (O, O)$	$\frac{c_s}{v_2-v_1} \geq R_d \geq \frac{c_s}{v_4-v_3}$	$P_A(t_N) < \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, P_A(t_N) \geq \frac{R_o \cdot v_a}{c_w},$ $\gamma < \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, \gamma \geq \frac{v_a \cdot R_o}{c_w}$
E14	$(m_H, m_H) - (O, O)$	$\frac{c_s}{v_4-v_3} > R_d > \frac{c_s}{v_2-v_1}$	
E15	$(m_N, m_N) - (O, R)$	$\frac{c_w}{v_1-v_2} \geq R_d \geq \frac{c_s}{v_4-v_3}$	$P_A(t_N) < \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, P_A(t_N) \geq \frac{R_o \cdot v_a}{c_w},$ $\gamma < \frac{v_a \cdot R_a}{c_c}, \gamma \geq \frac{v_a \cdot R_o}{c_w}$
E16	$(m_H, m_H) - (O, R)$	$\frac{c_s}{v_4-v_3} > R_d > \frac{c_w-c_s}{v_1-v_2}$	
E17	$(m_N, m_N) - (R, A)$	$\frac{c_s+c_c}{v_2-v_1} \geq R_d \geq \frac{c_s+b_w}{v_4-v_3}$	$P_A(t_N) < \frac{R_o \cdot v_a}{c_c}, P_A(t_N) < \frac{R_o \cdot v_a}{c_w},$ $\gamma \geq \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, \gamma \geq \frac{v_a \cdot R_a}{c_c}$
E18	$(m_H, m_H) - (R, A)$	$\frac{c_s+b_w}{v_4-v_3} > R_d > \frac{c_s+c_c}{v_2-v_1}$	
E19	$(m_N, m_N) - (R, O)$	$\frac{c_s+c_w-c_c}{v_2-v_1} \geq R_d \geq \frac{c_s}{v_4-v_3}$	$P_A(t_N) < \frac{R_o \cdot v_a}{c_c}, P_A(t_N) < \frac{R_o \cdot v_a}{c_w},$ $\gamma < \frac{v_a \cdot (R_a - R_o)}{c_c - c_w}, \gamma \geq \frac{v_a \cdot R_o}{c_w}$
E20	$(m_H, m_H) - (R, O)$	$\frac{c_s}{v_4-v_3} > R_d > \frac{c_s+c_w}{v_2-v_1}$	
E21	$(m_N, m_N) - (R, R)$	$\frac{c_s}{v_2-v_1} \geq R_d \geq \frac{c_c}{v_4-v_3}$	$P_A(t_N) < \frac{R_o \cdot v_a}{c_c}, P_A(t_N) < \frac{R_o \cdot v_a}{c_w},$ $\gamma < \frac{v_a \cdot R_a}{c_c}, \gamma < \frac{v_a \cdot R_o}{c_w}$
E22	$(m_H, m_H) - (R, R)$	$\frac{c_s}{v_4-v_3} > R_d > \frac{c_s}{v_2-v_1}$	

左の列から均衡の名前、防御者と攻撃者の戦略の組み合わせ、均衡条件、事前確率と事後確率の範囲である。 t_N と t_H の両タイプが m_H を送信する場合、表の γ は μ になる。

一括均衡は 16 種類存在し、4 種類の分離均衡と合わせて合計 22 種類の均衡が存在する。しかし中には現実的にあり得ない条件の均衡も複数存在するため、それらの均衡を省いて分析する必要がある。

5 結論と今後の課題

本論文では、ハニーポットを利用した防御者と DoS 攻撃者の間の相互作用を捉えた多段階シグナリングゲームを検討し、第1ステージの均衡解を解析した。本論文で提案したモデルを応用することにより、ネットワークセキュリティ分野において様々な欺瞞的防御を運用するための基本的な指針となると考えている。今後は多段階シグナリングゲームで数値実験を行い、欺瞞信号の減衰係数を含めた時にどのように戦略選択が変化していくのか分析したい。

参考文献

- [1] Basagiannis, S., P. Katsaros, A. Pombortsis, N. Alexiou (2009) “Probabilistic model checking for the quantification of DoS security threats”, *Computers and Security*, vol. 28, no. 6, pp. 450-465.
- [2] Carroll, T. E., D. Grosu (2009) “A Game Theoretic Investigation of Deception in Network Security”, *2009 Proceedings of 18th International Conference on Computer Communications and Networks*, pp. 1-6.
- [3] Ceker, H., J. Zhuang, S. Upadhyaya, Q. D. La, B. Soong (2016) “Deception-Based Game Theoretical Approach to Mitigame DoS Attacks”, *Lecture Notes in Computer Science*, vol. 9996, pp. 18-38.
- [4] Cho, I. K., D. M. Kreps (1987) ”Signaling games and stable equilibria”, *Quarterly Journal of Economics*, vol. 102, pp. 179-221.
- [5] Jiang, Z., Y. Ge, Y. Li (2005) “Max-utility wireless resource management for best-effort traffic”, *IEEE Transactions on Wireless Communications*, vol. 4, no. 1, pp. 100-111.
- [6] Meadows, C. (2001) “A cost-based framework for analysis of denial of service in networks”, *Journal of Computer Security*, vol. 9, no. 1, pp. 143-164.
- [7] Zhang, H. W., T. Li (2017) “Optimal active defense based on multi-stage attack-defense signaling game”, *Acta Electronica Sinica*, vol. 45, no. 2, pp. 431-439.