

# クラウド型 CAPTCHA に対する bot によるアクセスの分析

寺田 健太\* 岡部 寿男†  
京都大学

松本 悦宜‡  
Capy 株式会社

## 1 はじめに

近年、Web サイトに対する bot による不正アクセスが増加傾向にある。bot による不正アクセスの手口の典型として、リスト型攻撃によるパスワードクラッキングが挙げられる。また、不正アクセスとは言えないものの、転売目的で bot を用いてチケット買占めを行うような悪質なアクセスも見られるようになってきている。

bot による自動アクセスを防ぐ手段に CAPTCHA (Completely Automated Public Turing test to tell Computers and Humans Apart) がある。これはアクセスを行っているのが人なのか bot なのかを判別するためのテストである。しかし、bot の高度化に伴って、CAPTCHA の複雑化が必要となり、それに伴ってユーザの利便性が低下する問題がある。

そのため、リスクベース認証の考え方で、アクセス元の IP アドレスやブラウザの種類等の情報を基に bot であることが疑われるアクセスに対しては難しい CAPTCHA を課すことが行われている。[1]

本研究では、Capy 株式会社が提供するクラウド型 CAPTCHA サービスであるパズル CAPTCHA への実際のアクセスを分析し、リスクベース認証に有効な特徴量を検討した。ある日本国内向けの Web サイトに対するアクセスログについて、そのうち実際のインシデントから bot であると判断されたアクセスを分析したところ、海外の ISP からのアクセスで、かつブラウザが送出する HTTP Accept-Language が中国語であるという特徴が見つかり、それを手掛かりに他のアクセスについても精査したところ同様に疑わしいアクセスを多く見つけることができた。

## 2 Capy パズル CAPTCHA

Capy 社のパズル CAPTCHA はクラウド型のサービスで CAPTCHA に関するシステムは Web サイトとは独立のクラウド上のサーバに存在する。実装は Web サイトに CAPTCHA に関するスクリプトを追加すること

で行われる。

パズル CAPTCHA においてユーザによる通信は CAPTCHA に関するスクリプトの含まれるページのリクエストを行う get-js、パズル CAPTCHA で用いる画像のリクエストを行う get-image、パズル CAPTCHA の回答のポストを行う verify の3手順からなり、Web サイトや CAPTCHA 生成/照合サーバと通信を行う。[2]

## 3 データセット

今回分析に用いたデータセットは、ある顧客のサービス上におけるパズル CAPTCHA へのアクセスログ 622077 件に bot によるアクセスかどうかを判断したフラグを付与したものである。この判断は Capy 社が通常取得しているログに加えて、導入元自身による調査も含め、Capy 社により行われている。なおアクセス期間は 2019 年 3 月 1 日 9 時 0 分 0 秒から 2019 年 4 月 1 日 8 時 59 分 59 秒までである。このサービスは日本国内向けのもので、bot による不正なアクセスによって金銭的被害が想定されるものである。

データセットに含まれる情報のうち、本研究ではアクセス時刻、アクセス元の IP アドレス、HTTP User-Agent、HTTP Accept-Language、パズル認証の成否、パズル認証の回答時間、ならびに is\_bot フラグを用いた。is\_bot フラグは先述した bot によるアクセスであると判断されたアクセスログに立てられるフラグである。

またアクセス件数やそのグラフについては get-image のアクセスのみをカウントすることにする。

## 4 データの分析

### 4.1 Bot と判断されたアクセスの分析

まず is\_bot フラグが True のアクセスログ (426 件) について、HTTP Accept-Language を分析したところ、HTTP Accept-Language が “zh-CN,zh;q=0.9” または “zh-CN,zh;q=0.8” のいずれかに限られることが判明した “zh” は中国語、“zh-CN” は中国語 (中国本土) を表す言語コードである。また IP アドレスについて whois サービスで調べた結果、香港にある特定の ISP に割り当てられているアドレスからのアクセスであることがわかった。これは国内向けの日本語でのサービスへのアク

Analyzing bot access to cloud-based CAPTCHA

\* KENTA TERADA, Kyoto University

† YASUO OKABE, Kyoto University

‡ YOSHINORI MATSUMOTO, Capy Japan Inc.

セスとしては極めて不審なアクセスであると言える。

また、これらのアクセスの時期は、分析の対象としている期間において、全体のアクセスの分布と比較すると、3月下旬に集中していることが言える(図1、図2のbot)。

一方、それ以外の情報については顕著な特徴は見出せなかった。

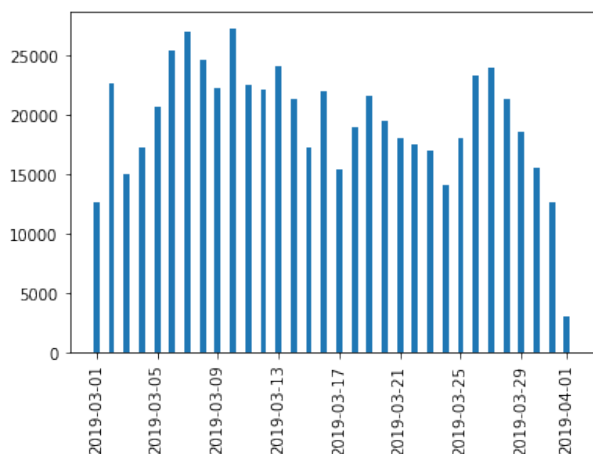


図1: 全体のアクセスの分布

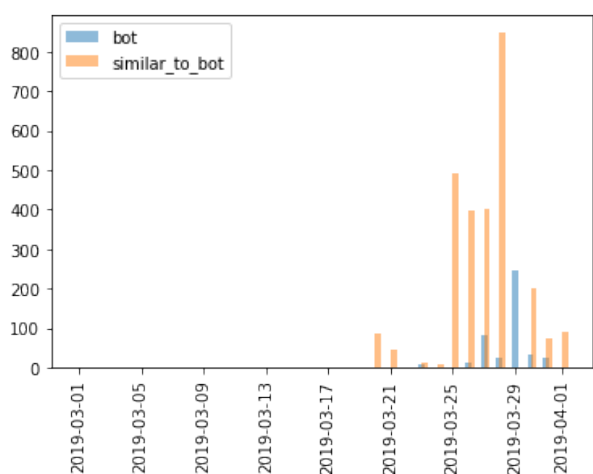


図2: bot とそれに類似したアクセスの分布

#### 4.2 Bot に類似したアクセスの分析

次に is\_bot フラグが False であるアクセスのうち、HTTP Accept-Language が is\_bot フラグが True で見られたものと同じであるアクセス (7693 件) の IP アドレスを whois サービスで調べた。その結果、is\_bot フラグが True のものと同一の ISP から割り当てられている IP アドレスが見つかった。それらの IP アドレスと一致するアクセスを is\_bot フラグが False のものから抽出したところ 2649 件見付き、HTTP Accept-Language

は is\_bot フラグが True のアクセスと同様、すべて “zh-CN,zh;q=0.9” または “zh-CN,zh;q=0.8” のいずれかに限られていた。これらは is\_bot フラグが True のアクセスと同程度に bot によるアクセスであることが疑われ、かつ同一の攻撃者グループによる同一の ISP からの攻撃である可能性が高い。またアクセス時期は is\_bot が True であるアクセスに少し先立つ形で分布しており(図2の similar\_to\_bot)、ログ調査により bot と判断されたアクセスよりも前から攻撃が始まっていたと推定されている。

さらに、HTTP Accept-Language が、is\_bot が True のものと同じで IP アドレスが上記の ISP 以外の範囲にあるアクセス (5044 件) についての分析を行ったところ、その大半 (4899 件) が中国にある特定の ISP に割り当てられたアドレスからのアクセスであり、それ以外は香港の ISP から 122 件、国内の ISP のアドレスから 23 件と比較的少なかった。

## 5 まとめ

本研究では、パズル CAPTCHA への bot によるアクセスの実際のログを分析し、Web サイトの対象の国や言語と、アクセス元の IP アドレスが示す地域ならびにブラウザの受け入れる言語の不一致が bot の判別に有効であることを示した。また、それを用いてさらに bot によるアクセスと思われる他のアクセスを発見することができた。

今後は更に多言語かつ全世界向けの Web サイトについて同様の分析を進める予定である。

## 参考文献

- [1] S. Wiefeling, L. Lo Iacono, M. Dürmuth, “Is This Really You? An Empirical Study on Risk-Based Authentication Applied in the Wild,” In: Dhillon G., Karlsson F., Hedström K., Zúquete A. (eds) ICT Systems Security and Privacy Protection. SEC 2019. IFIP Advances in Information and Communication Technology, vol 562. Springer, Cham. [https://doi.org/10.1007/978-3-030-22312-0\\_10](https://doi.org/10.1007/978-3-030-22312-0_10)
- [2] T. Arai, Y. Okabe, Y. Matsumoto and K. Kawamura, “Detection of Bots in CAPTCHA as a Cloud Service Utilizing Machine Learning,” 2020 International Conference on Information Networking (ICOIN), Barcelona, Spain, 2020, pp. 584-589, doi: 10.1109/ICOIN48656.2020.9016522.