

Tor ネットワークのクローラシステム (2) CAPTCHA 自動解析

梶間 大地 †

鳥居 洸希 †

菊池 浩明 †

明治大学総合数理学部 †

1 はじめに

近年、セキュリティへの意識から匿名通信システム Tor(The Onion Routing) の利用が増えている。Tor は匿名での通信が可能のため違法商品販売等の不正な目的で利用されている場合も多く、これらの違法商品販売サイトの調査を試みている。しかしダークウェブの違法商品販売サイトの多くは CAPTCHA を採用しているため機械的な調査を行うことができなかった [1]。

そこで本研究では、OCR を用いて Tor ネットワーク内に存在する違法商品販売サイトを調査を行う。本稿は、OCR を用いた CAPTCHA 画像の自動解析システムの開発について報告する。システム全体と人間が CAPTCHA と解く操作のみを行う半自動クローラについては [2] で報告する。ただし、本研究で扱う CAPTCHA の種類はテキストベースのものに限る。

2 全自動クローラシステムの開発

本研究では、Tor ネットワークにアクセスするために、Tor をインストールした Linux サーバと Python を用いた。図 1 に本研究で開発したクローラシステムの構成図を示す。

Python プログラムに対して Tor ネットワーク内のク

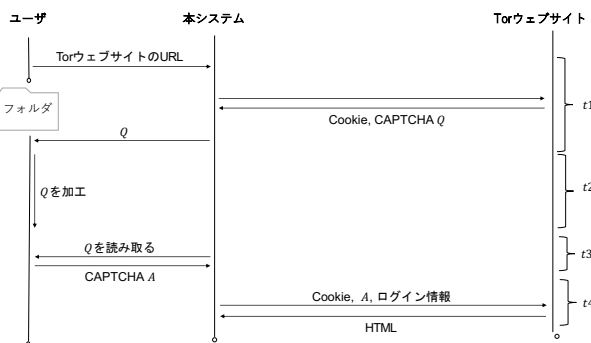


図1 システム構成図

ローリングを行いたいサイトの URL を送る。URL を受け取った Python プログラムは socks5, Tor を利用してサイトにアクセスする。アクセスした際に表示された CAPTCHA 画像 Q をサーバ内に保存し、セッションの Cookie 情報を取得する。次にサーバに保存された Q を Python プログラムを用いて画像加工を行い、加工した Q をサーバに保存する。加工した Q を Tesseract OCR を利用して文字認識を行い、帰ってきた文字列 A を Python プログラムに渡す。その後、サイトにユーザー名、パスワードのログイン情報、保存した Cookie, A を利用してログインを試みる。サイトのユーザーページにアクセスできれば終了。 A が誤りの場合、プログラムに URL を送った時点からやり直す。これをユーザーのページにアクセスできるまで繰り返す。

3 CAPTCHA の機械解読

3.1 CAPTCHA の成功率

表 1 に CAPTCHA の解決の成功率を示す。

3.2 CAPTCHA 画像の加工手法

CAPTCHA 画像の解析には openCV の Python 用ライブラリを用いて画像加工を行い、google 社が開発を行った OCR(光学的文字認識) ツールである Tesseract OCR[3] を用いて加工後の画像の文字認識を行った。画像の加工処理では、取得した CAPTCHA 画像を 2 値化処理を行ったのち、閾値処理によってノイズの除去を行った。また、CAPTCHA 画像にノイズのないものには閾値処理を行わず 2 値化処理のみとした。



図2 加工後の CAPTCHA 画像の例

表1 CAPTCHA 解読の成功率







URL	主要カテゴリ	全体	正答	誤答	認識なし	閾値	CAPTCHA 画像の例
apollonujscjrIng.onion	Credit Card	1000	0	722	278	120	
deluxedzn6h572qp.onion	内容不明	1000	17	983	0	120	
27kaqicipyhous2p.onion	Credit Card	1000	653	347	0	なし	
abyssopyps3z4xof.onion	Drugs	1000	0	467	533	120	
yohoho7do5u6d35g.onion	Drugs	1000	303	697	0	なし	
torevvq44o7ofjuu.onion	Credit Card	1000	596	331	73	120	

表2 HTML 取得にかかった時間

	平均 [秒]	標準偏差 [秒]
t1	2.717	0.286
t2	0.013	0.002
t3	0.195	0.030
t4	2.444	0.187
全体	7.237	3.483

3.3 処理時間の評価

表2にHTML取得にかかった時間を計測した結果を示す。これはtorevvq44o7ofjuu.onionに合計150回アクセスした際のデータである。

3.4 考察

自動解読できなかったCAPTCHA画像については以下のような原因が考えられる。

1. 図3のように読み込む文字とノイズの色が似ていたりノイズが文字の上に存在しているため、閾値処理を行った際に読み解く文字まで欠けたり削除されてしまう。

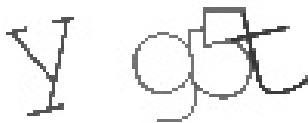


図3 考察1の例

2. 読み解く文字の歪みや重なりによって、文字認識に失敗する

4 おわりに

本研究では、Torネットワーク内に存在する違法商品販売サイトの情報を効率的に収集するための全自動クロールシステムの開発を行った。CAPTCHA画像については平均26.15%で解決できることがわかった。しかし、汎用的に用いることは出来ないこともわかった。今後はCAPTCHA画像を加工すること以外でのCAPTCHAの自動解決の方法を考案することを課題とする。

参考文献

- [1] 鳥居洗希, 菊池浩明. “Torネットワークのクロールシステムの開発と違法商品販売サイトの調査”, 情報処理学会第81回全国大会 pp.3.435-3.436, 2019.
- [2] 鳥居洗希, 菊池浩明. “Torネットワークのクロールシステム(1)違法商品販売の調査”, 情報処理学会第82回発表予定.
- [3] Tesseract OCR, <https://opensource.google/projects/tesseract> (2019年12月参照)