# 윈도우 환경에서의 **웹크롤링을 통한 악성코드 검사**

## 바이러스토탈 API DB연동과 광고성 웹크롤링

**SAMZO**

1771064 김연희 1771070 안해인
1771071 유예린 1771077 전재원

SAMZO

# 목차

# 1. 프로젝트 소개

# [구현목적]

# [개발환경] 1. PyCharm



- 파이썬으로 개발할 때 사용하는 대표적인 도구
- 편리한 인터페이스제공

→ 웹크롤링에 사용

# [개발환경] 2. Virustotal



-   의심스런 파일과 URL을 분석하고
바이러스, 웜, 트로얀과 모든 종류의 악성 코드를 쉽고,
빠르게 탐지할 수 있는 편리한 무료 서비스
- 여러 개의 백신 엔진으로 검사하여 그 결과를 투명하게
보여줌

→Virustotal API를 사용하여 크롤링한 URL 검사

# [개발환경] 3. Elasticsearch

SAMZO

elasticsearch

- 아파치 루신을 기반으로 만든 <u>분산 검색엔진(오픈소스)</u>
- 분산 시스템이기 때문에 검색 대상 용량이
증가했을 때 대처하기 쉬움
- 깃허브, 이베이, 가디언 같은 기업이
엘라스틱서치 기술로 내부 검색 기능을 구축함

- <u>데이터를 저장</u>할 수 있어 NoSQL 저장소로도 활용할 수 있음
- <u>웹 서버나 데이터베이스 서버에 수정&삭제를</u>
<u>요청하고 결과를 받을 수 있는 응용 프로그래밍</u>
<u>인터페이스(API)를 지원함</u>

→ Virustotal 분석 결과 저장&검색

# [개발환경] 4. Kibana

SAMZO

Kibana

- 엘라스틱 서치의 대시보드 어플리케이션
- Elasticsearch 데이터의 시각화와 Elastic Stack의 탐색을 지원

SAMZO

# 2. 구현

# [웹크롤링]

# [웹크롤링]

# [웹크롤링]

# [웹크롤링]

# [웹크롤링]

SAMZO

[윈도우 환경에서 엘라스틱서치 DB 사용]

# [윈도우 환경에서 엘라스틱서치 DB 사용]

Elasticsearch를 다운 받은 후, elasticsearch.bat을 실행시킨다.

# [윈도우 환경에서 엘라스틱서치 DB 사용]

Python에서 elasticsearch를 사용하기 위해 pip install을 한 후, 인덱스 생성

SAMZO

# [바이러스토탈]

# [바이러스토탈]

# [바이러스토탈]

해당 URL을 바이러스토탈로 분석 후 나온 결과값을 출력

```python
import requests
headers = {
    "Accept-Encoding": "gzip, deflate",
    "User-Agent": "gzip,  My Python requests library example client or username"
}
params = {'apikey': '1cfba73f4664a297c8c16a3a358fad56b0c4997a0ccf06dad74a44cca3ba4d11', 'resource':'http://alicdn.re.kr/re.kr'}
response = requests.post('https://www.virustotal.com/vtapi/v2/url/report', params=params, headers=headers)
json_response = response.json()
j = json_response

print (j)
```

==================== RESTART: C:\Users\안해인\Desktop\확인용.py ====================
{'scan_id': 'f9f6233c1840eb98c87f7df7fc85927dd4cc0f31faac32b7a960ef5cc8c66f61-1514263742', 'resource': 'http://alicdn.re.kr/re.kr', 'url': 'http://alicdn.re.kr/re.kr', 'response_code': 1, 'scan_date': '2017-12-26 04:49:02', 'permalink': 'https://www.virustotal.com/url/f9f6233c1840eb98c87f7df7fc85927dd4cc0f31faac32b7a960ef5cc8c66f61/analysis/1514263742/', 'verbose_msg': 'Scan finished, scan information embedded in this object', 'filescan_id': None, 'positives': 1, 'total': 66, 'scans': {'CLEAN MX': {'detected': False, 'result': 'clean site'}, 'DNS8': {'detected': False, 'result': 'clean site'}, 'OpenPhish': {'detected': False, 'result': 'clean site'}, 'VX Vault': {'detected': False, 'result': 'clean site'}, 'ZDB Zeus': {'detected': False, 'result': 'clean site'}, 'ZCloudsec': {'detected': False, 'result': 'clean site'}, 'PhishLabs': {'detected': False, 'result': 'unrated site'}, 'Zerofox': {'detected': False, 'result': 'clean site'}, 'K7AntiVirus': {'detected': False, 'result': 'clean site'}, 'FraudSense': {'detected': False, 'result': 'clean site'}, 'Virusdie External Site Scan': {'detected': False, 'result': 'clean site'}, 'Quttera': {'detected': False, 'result': 'clean site'}, 'AegisLab WebGuard': {'detected': False, 'result': 'clean site'}, 'MalwareDomainList': {'detected': False, 'result': 'clean site', 'detail': 'http://www.malwaredomainlist.com/mdl.php?search=alicdn.re.kr'}, 'ZeusTracker': {'detected': False, 'result': 'clean site', 'detail': 'https://zeustracker.abuse.ch/monitor.php?host=alicdn.re.kr'}, 'zvelo': {'detected': False, 'result': 'clean site'}, 'Google Safebrowsing': {'detected': False, 'result': 'clean site'}, 'Kaspersky': {'detected': False, 'result': 'unrated site'}, 'BitDefender': {'detected': False, 'result': 'clean site'}, 'Opera': {'detected': False, 'result': 'clean site'}, 'Certly': {'detected': False, 'result': 'clean site'}, 'G-Data': {'detected': False, 'result': 'clean site'}, 'C-SIRT': {'detected': False, 'result': 'clean site'}, 'CyberCrime': {'detected': False, 'result': 'clean site'}, 'SecureBrain': {'detected': False, 'result': 'clean site'}, 'Malware Domain Blocklist': {'detected': False, 'result': 'clean site'}, 'MalwarePatrol': {'detected': False, 'result': 'clean site'}, 'Webutation': {'detected': False, 'result': 'clean site'}, 'Trustwave': {'detected': False, 'result': 'clean site'}, 'Web Security Guard': {'detected': False, 'result': 'clean site'}, 'CyRadar': {'detected': False, 'result': 'clean site'}, 'desenmascara.me': {'detected': False, 'result': 'clean site'}, 'ADMINUSLabs': {'detected': False, 'result': 'clean site'}, 'Malwarebytes hpHosts': {'detected': False, 'result': 'clean site'}, 'Dr.Web': {'detected': False, 'result': 'clean site'}, 'AlienVault': {'detected': False, 'result': 'clean site'}, 'Emsisoft': {'detected': False, 'result': 'clean site'}, 'Rising': {'detected': False, 'result': 'clean site'}, 'MalcOde Database': {'detected': False, 'result': 'clean site', 'detail': 'http://malcOde.com/database/index.php?search=alicdn.re.kr'}, 'malwares.com URL checker': {'detected': False, 'result': 'clean site'}, 'Phishtank': {'detected': False, 'result': 'clean site'}, 'Malwared': {'detected': False, 'result': 'clean site'}, 'StopBadware': {'detected': False, 'result': 'unrated site'}, 'Antiy-AVL': {'detected': False, 'result': 'clean site'}, 'Forcepoint ThreatSeeker': {'detected': False, 'result': 'clean site'}, 'SCUMWARE.org': {'detected': False, 'result': 'clean site'}, 'Comodo Site Inspector': {'detected': False, 'result': 'clean site'}, 'Malekal': {'detected': False, 'result': 'clean site'}, 'ESET': {'detected': False, 'result': 'clean site'}, 'Sophos': {'detected': False, 'result': 'unrated site'}, 'Yandex Safebrowsing': {'detected': False, 'result': 'clean site', 'detail': 'http://yandex.com/infected?l10n=en&url=http://alicdn.re.kr/re.kr'}, 'Spam404': {'detected': False, 'result': 'clean site'}, 'Nucleon': {'detected': False, 'result': 'clean site'}, 'Sucuri SiteCheck': {'detected': True, 'result': 'malicious site'}, 'Blueliv': {'detected': False, 'result': 'clean site'}, 'Netcraft': {'detected': False, 'result': 'unrated site'}, 'AutoShun': {'detected': False, 'result': 'unrated site'}, 'ThreatHive': {'detected': False, 'result': 'clean site'}, 'FraudScore': {'detected': False, 'result': 'clean site'}, 'Tencent': {'detected': False, 'result': 'clean site'}, 'URLQuery': {'detected': False, 'result': 'unrated site'}, 'Fortinet': {'detected': False, 'result': 'clean site'}, 'ZeroCERT': {'detected': False, 'result': 'clean site'}, 'Baidu-International': {'detected': False, 'result': 'clean site'}, 'securolytics': {'detected': False, 'result': 'clean site'}}}

# [바이러스토탈]

바이러스토탈 분석에서 malware가 뜨는 결과만 추출

```python
import requests
headers = {
    "Accept-Encoding": "gzip, deflate",
    "User-Agent": "gzip,  My Python requests library example client or username"
}
params = {'apikey': '1cfba73f4664a297c8c16a3a358fad56b0c4997a0ccf06dad74a44cca3ba4d11', 'resource':'http://alicdn.re.kr/re.kr'}
response = requests.post('https://www.virustotal.com/vtapi/v2/url/report', params=params, headers=headers)
json_response = response.json()
j = json_response

if(j['response_code']!=0):
        scan_id=j['scan_id']
        url=j['resource']
        result=j['scans']
        tester=list(result.keys())
        val=list(result.values())
        result=[d['detected']for d in val]
        mal=[d['result']for d in val]

        print("\n\nscan_id: ", scan_id)
        print("URL: ",url)
        print("\n==========Detection==========")
        for i in range(0,65):
                if result[i] == True:
                    print (tester[i],result[i], mal[i],"\n")

        i=i+1
```

```
scan_id:   f9f6233c1840eb98c87f7df7fc85927dd4cc0f31faac32b7a960ef5cc8c66f61-1514263742
URL:  http://alicdn.re.kr/re.kr

==========Detection==========
Sucuri SiteCheck True malicious site
```

# [바이러스토탈]

```python
import re
import requests
from bs4 import import BeautifulSoup
import json
from elasticsearch import Elasticsearch
es= Elasticsearch()
es.indices.create(index='vv',body='')

headers = {
    "Accept-Encoding": "gzip, deflate",
    "User-Agent": "gzip,  My Python requests library example client or username"
}

siteurl = "http://news1.kr/articles/?3488307"
html = requests.get(siteurl).text
soup = BeautifulSoup(html ,"html.parser")
adbox=soup.find("div",{"class":"ad_wing_l"})

for src in adbox.select('div.wingLeftBlock a[href]'):|
    adurl=src['href']
    params = {'apikey': '1cfba73f4664a297c8c16a3a358fad56b0c4997a0ccf06dad74a44cca3ba4d11', 'resource': adurl}
    response = requests.post('https://www.virustotal.com/vtapi/v2/url/report', params=params, headers=headers)
    json_response = response.json()
    j = json_response

    if(j['response_code']!=0):
        scan_id=j['scan_id']
        url=j['resource']
        result=j['scans']
        tester=list(result.keys())
        val=list(result.values())
        result=[d['detected']for d in val]
        mal=[d['result']for d in val]

        print("\n\nscan_id: ", scan_id)
        print("URL: ",url)
        print("\n==========Detection==========")
        for i in range(0,65):
            if result[i] == True:
                print (tester[i],result[i], mal[i],"\n")
            i=i+1
for k in range(0,65):
    if result[k] == True:
        doc = {'sitename' : scan_id, 'url' : url, 'tester' : tester[k], 'result' : result[k], 'malware' : mal[k]}
    else:
        doc = {'sitename' : scan_id, 'url' : url, 'tester' : tester[k], 'result' : result[k], 'malware' : mal[k]}
    k=k+1
    es.index(index="vv", doc_type="virustotal",body=doc)
```

```
scan_id:  d5deb3116b10e59e425d70924c848b16e2eba40128aed3237de3c2b429d275bb-1539761721
URL:  https://www.cryptohub.or.kr

==========Detection==========


scan_id:  8ffd1b716debb45961d20d0f3a94d02d1afb20138c483d0d2ee26c5b482937c7-1525960952
URL:  http://www.happypet.co.kr/

==========Detection==========


scan_id:  1b276d8cab9d9a7436af1b488078e13f51974adec3612493f1cee453e1a8188f-1527160392
URL:  http://ncook.news1.kr/

==========Detection==========
>>>
```

# [바이러스토탈]

## Elasticsearch DB에 삽입한 데이터를 kibana로 연동시킴

sitename: d5deb3116b10e59e425d70924c848b16e2eba40128aed3237de3c2b429d275bb-1539761721  url: https://www.cryptohub.or.kr  tester: CLEAN MX, DNS8, OpenPhish, VX Vault, ZDB Zeus, ZCloudsec, PhishLabs, Zerofox, K7AntiVirus, FraudSense, Virusdie External Site Scan, Quttera, AegisLab Web Guard, MalwareDomainList, ZeusTracker, zvelo, Google Safebrowsing, Kaspersky, BitDefender, Opera, Certly, G-Data, C-SIRT, CyberCrime, SecureBrain, Malware Domain Blocklist, MalwarePatrol, Webutation, Trustwave, Web Security Guard, CyRadar, desenmascara.me, ADMINUSLabs, Malwarebytes hpHosts, Dr.Web, AlienVault, Emsisoft, Rising, Malc0de Database, malwares.com URL checker, Phishtank, Malwared, Avira, NotMining, StopBadwar

Table | JSON

```json
1 {
2     "_index": "samz02",
3     "_type": "malware",
4     "_id": "YhIsnmcB5ib-0eM6eK0y",
5     "_version": 1,
6     "_score": 1,
7     "_source": {
8         "sitename": "d5deb3116b10e59e425d70924c848b16e2eba40128aed3237de3c2b429d275bb-1539761721",
9         "url": "https://www.cryptohub.or.kr",
10        "tester": [
11            "CLEAN MX",
12            "DNS8",
13            "OpenPhish",
14            "VX Vault",
15            "ZDB Zeus",
16            "ZCloudsec",
17            "PhishLabs",
18            "Zerofox",
19            "K7AntiVirus",
20            "FraudSense",
21            "Virusdie External Site Scan",
22            "Quttera",
23            "AegisLab WebGuard",
24            "MalwareDomainList",
25            "ZeusTracker",
26            "zvelo",
27            "Google Safebrowsing",
28            "Kaspersky",
29            "BitDefender",
30            "Opera",
31            "Certly",
32            "G-Data",
33            "C-SIRT",
34            "CyberCrime",
```

_source

sitename: 1b276d8cab9d9a7436af1b488078e13f51974adec3612493f1cee453e1a8188f-1527160392  url: http://ncook.news1.kr/  tester: DNS8  result: false  malware: clean site  _id: ORJHoWcB5ib-0eM6B645  _type: virustotal  _index: vvv  _score: 1

sitename: 1b276d8cab9d9a7436af1b488078e13f51974adec3612493f1cee453e1a8188f-1527160392  url: http://ncook.news1.kr/  tester: OpenPhish  result: false  malware: clean site  _id: OhJHoWcB5ib-0eM6B65Q  _type: virustotal  _index: vvv  _score: 1

Table | JSON

View single document

```json
1 {
2     "_index": "vvv",
3     "_type": "virustotal",
4     "_id": "OhJHoWcB5ib-0eM6B65Q",
5     "_version": 1,
6     "_score": 1,
7     "_source": {
8         "sitename": "1b276d8cab9d9a7436af1b488078e13f51974adec3612493f1cee453e1a8188f-1527160392",
9         "url": "http://ncook.news1.kr/",
10        "tester": "OpenPhish",
11        "result": false,
12        "malware": "clean site"
13    }
14 }
```

sitename: 1b276d8cab9d9a7436af1b488078e13f51974adec3612493f1cee453e1a8188f-1527160392  url: http://ncook.news1.kr/  tester: Virusdie External Site Scan  result: false  malware: clean site  _id: QhJHoWcB5ib-0eM6CK4U  _type: virustotal  _index: vvv  _score: 1

sitename: 1b276d8cab9d9a7436af1b488078e13f51974adec3612493f1cee453e1a8188f-1527160392  url: http://ncook.news1.kr/  tester: MalwareDomainList  result: false  malware: clean site  _id: RRJHoWcB5ib-0eM6CK5E  _type: virustotal  _index: vvv  _score: 1

sitename: 1b276d8cab9d9a7436af1b488078e13f51974adec3612493f1cee453e1a8188f-1527160392  url: http://ncook.news1.kr/  tester: zvelo  result: fal

# 3. 한계점&어려웠던 점

SAMZO

# [한계점]

✓ 사이트마다 html 태그나 클래스명, 구조 등 달라서 일반화할 수가 없었다.
  • 다른 사이트를 돌릴 때 마다 코드의 수정이 필요했고, 크롤링할 때 마다 힘들게 찾아야했다.

✓ 실제 사이트 광고에는 악성코드가 발견되지 않아 실질적으로 구현하고자 하는 결과를 확인할 수가 없다.
  • 악성URL만 따로 바이러스토탈에 돌려 실제 악성코드 발견시에는 어떻게 결과가 나오는지 부분적으로 확인이 가능하다.

# [어려웠던 점]

---

✓ 윈도우 환경에서의 활용이 어디에도 나와있지 않아 창조적으로 시도하게 되었다.

✓ 엘라스틱서치에 관련된 내용들이 리눅스기반 사용이어 서 윈도우 버전에 맞게 새로 코드를 작성해야 했다.

✓ 바이러스토탈 결과값에서 필요한 정보를 뽑기 위해 어떤 형식으로 출력되는지 실제로 코드를 짜면서 확인했어야 했다.

# [참고자료]

SAMZO

+ Google Search
+ 교수님

# SAMZO

---

# 감사합니다