

Статья Несколько слов об энтропии

 xss.is/threads/54022



Всем известно что энтропия - это мера хаоса системы, а "информация" (данные) - это мера упорядоченности системы. Используя это свойство можно выполнять качественный анализ данных вне зависимости от их характера, вычисляя функцию частоты распределения: значение этой функции будет принимать большие значения для данных содержащих больше "хаоса" и, соответственно, меньшие значения для данных, содержащих меньше "хаоса".

Немного теории

Программа, демонстрирующая описанное выше свойство частоты распределения байтов данных в файле приведена далее (оригинал лежащего в основе кода доступен [здесь](#)).

Python:

```

import sys
import math
import os
import numpy as np
import matplotlib.pyplot as plt

if len(sys.argv) != 2:
    print("Usage: " + os.path.basename(__file__) + " filename")
    sys.exit()

f = open(sys.argv[1], "rb")
byteArr = map(ord, f.read())
f.close()
fileSize = len(byteArr)
#print("File size in bytes: " + fileSize)

freqList = []
for b in range(256):
    ctr = 0
    for byte in byteArr:
        if byte == b:
            ctr += 1
    freqList.append(float(ctr) / fileSize)
#print("Frequencies of each byte-character: " + freqList)

# entropy
ent = 0.0
for freq in freqList:
    if freq > 0:
        ent = ent + freq * math.log(freq, 2)
ent = -ent

N = len(freqList)

ind = np.arange(N) # the x locations for the groups
width = 1.00      # the width of the bars

#fig = plt.figure()
fig = plt.figure(figsize=(11,5),dpi=100)
ax = fig.add_subplot(111)
rects1 = ax.bar(ind, freqList, width)
ax.set_autoscalex_on(False)
ax.set_xlim([0,255])

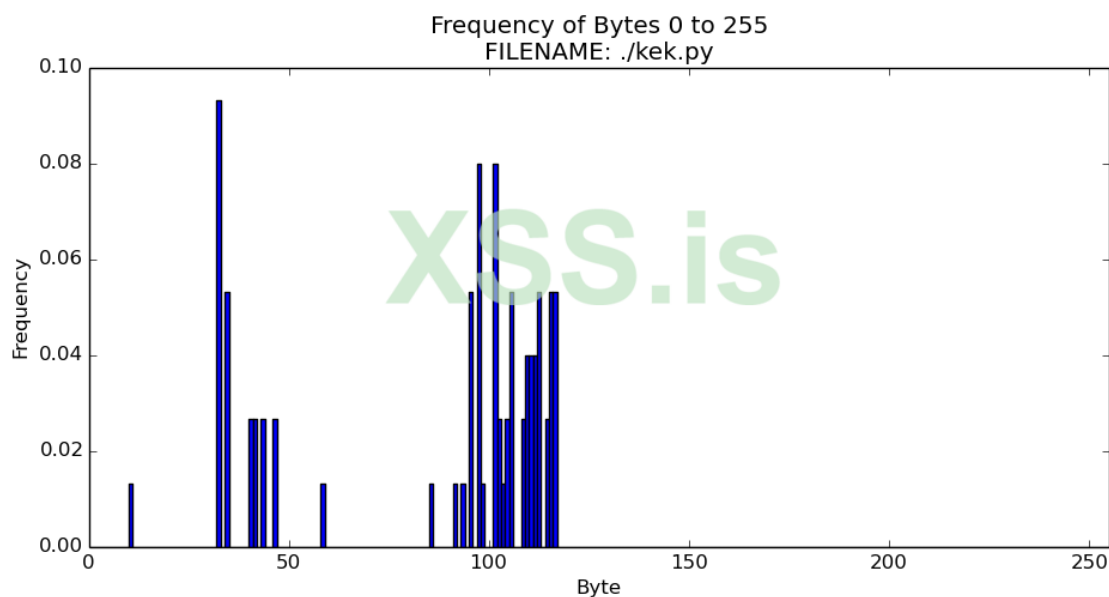
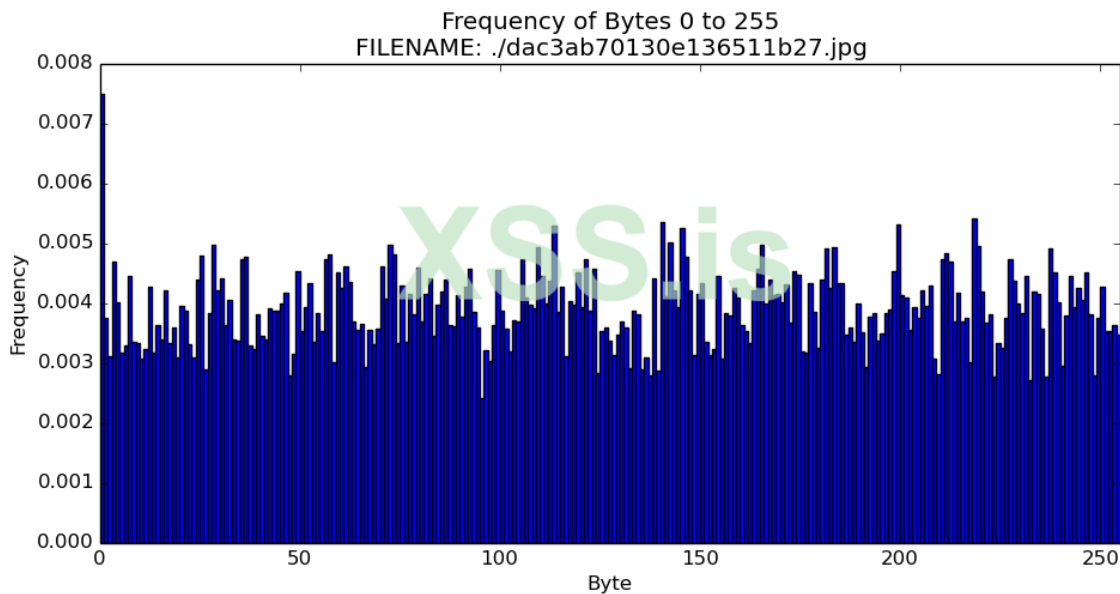
ax.set_ylabel('Frequency')
ax.set_xlabel('Byte')
ax.set_title('Frequency of Bytes 0 to 255\nFILENAME: ' + sys.argv[1])

plt.show()

```

В результате выполнения программы, для заданного файла, вычисляется значение энтропии, а также строится частотная характеристика.

В качестве наглядного примера можно рассмотреть частотные характеристики файла со сжатием и без сжатия (первый и второй графики соответственно).



Из графиков видно, что у сжатого файла энтропия выше, чем у файла без сжатия. Этот принцип широко *используется при анализе бинарных файлов, позволяя сделать вывод о наличии/отсутствии сжатия в файле* (и его типе).

Помимо анализа бинарных файлов свойство энтропии используется и при анализе обычных файлов: у строк, содержащих секреты (пароли/криптографические ключи/ токены и т.п.) энтропия выше чем у обычных строк. Применение энтропии при анализе данных существенно ускоряет этот процесс.

На этом краткий экскурс в теории можно закончить и перейти к практике.

Софт

Софта подобной тематики существует масса (любят велосипеды делать).

Наиболее известным считается truffleHog, в котором реализованы типовые функции для анализа (поиск по заданным шаблонам, поиск по энтропии и т.д.). Имеет большой недостаток - работает только с git-репозиториями, без возможности анализа локальных файлов. Существует доработанная до ума модификация - truffleHog3. Среди доработок: добавили возможность работы с локальными файлами, а также функции фильтрации результатов. Как и оригинал, truffleHog3 также давно не обновлялся.

В настоящее время активно развивается идейный продолжатель truffleHog - rustyHog, превосходящий предшественника как по скорости работы, так и по количеству реализованных фич.

Распространяется rustyHog в виде исходных текстов (компилируется в несколько команд, как под windows, так и под linux/macos), также имеется готовый образ в Docker.

Для удобства использования rustyHog разбит на компоненты по области применения:

- Ankamali Hog: для работы с гугл-дискон.
- Berkshire Hog: для работы с S3-инстансами.
- Choctaw Hog: для работы с Git-репозиториями.
- Duroc Hog: для работы с файловой системой.
- Essex Hog: для работы с Confluence.
- Gottingen Hog: для работы с JIRA.

Список параметров у каждой из утилит типовой, пример справки для duroc_hog приведён на рисунке ниже.



```
resultCd -- -zsh -- 160x50
~/Desktop/resultCd -- -zsh
Last login: Wed Jul 14 23:16:11 on ttys003
[estri@MacBook-Pro-estri resultCd % duroc_hog -h
duroc_hog 1.0.10
Scott Cutler <scutler@newrelic.com>
File system secret scanner in Rust

USAGE:
  duroc_hog [FLAGS] [OPTIONS] <FSPATH>

FLAGS:
  --caseinsensitive  Sets the case insensitive flag for all regexes
  --entropy          Enables entropy scanning
  --norecursive     Disable recursive scanning of all subdirectories underneath the supplied path
  --prettyprint     Outputs the JSON in human readable format
  -z, --unzip       Recursively scans archives (ZIP and TAR) in memory (dangerous)
  -v, --verbose     Sets the level of debugging information
  -h, --help       Prints help information
  -V, --version    Prints version information

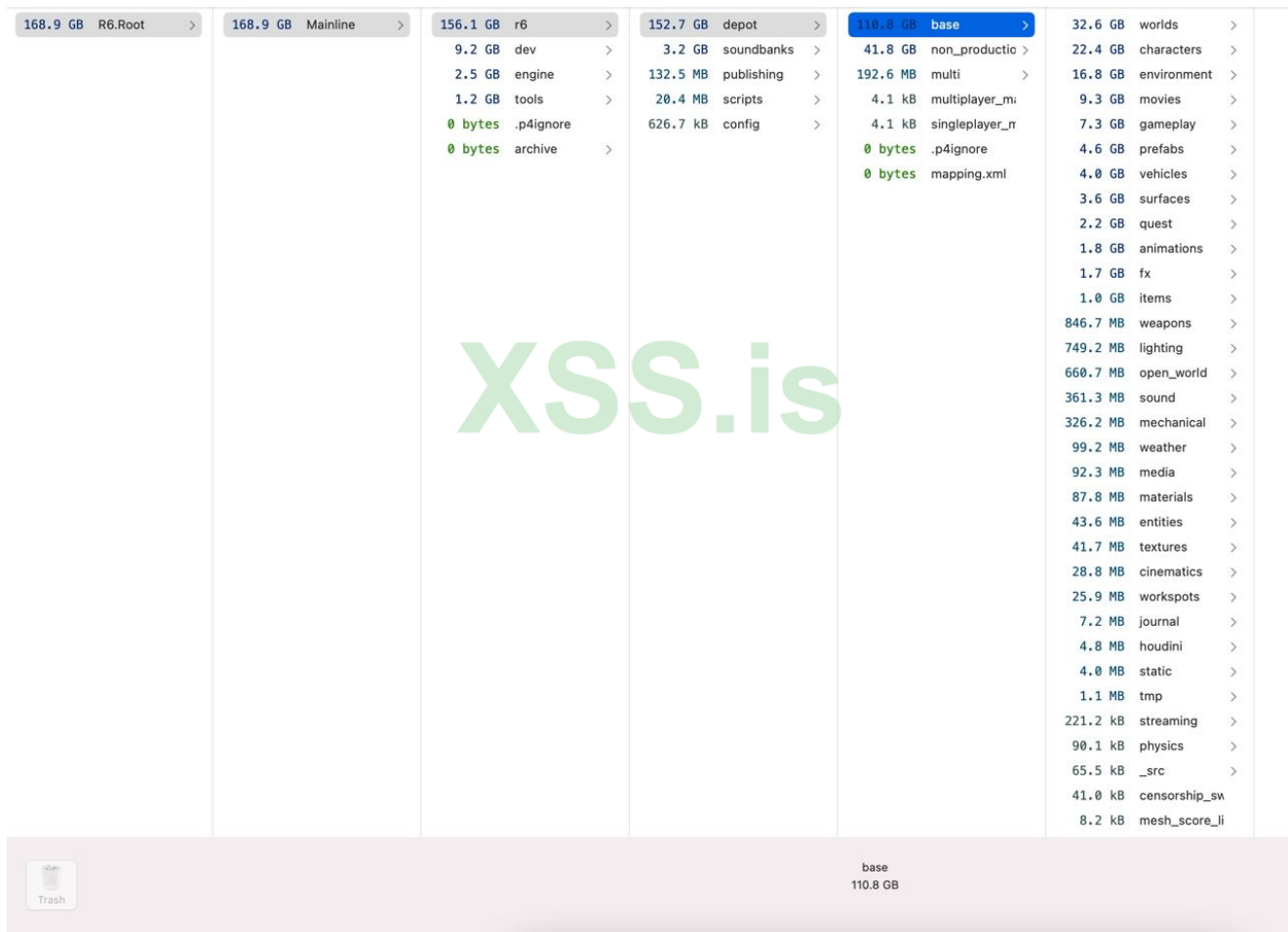
OPTIONS:
  -a, --allowlist <ALLOWLIST>      Sets a custom allowlist JSON file
  --default_entropy_threshold <DEFAULT_ENTROPY_THRESHOLD>  Default entropy threshold (0.6 by default)
  -o, --outputfile <OUTPUT>        Sets the path to write the scanner results to (stdout by default)
  -r, --regex <REGEX>              Sets a custom regex JSON file

ARGS:
  <FSPATH>  Sets the path of the directory or file to scan.
estri@MacBook-Pro-estri resultCd %
```

Анализ данных

Чтобы продемонстрировать возможности анализа необходим хороший и наглядный пример: слитая в паблик утечка CD-Project отлично для этого подойдёт.

В числе выложенных и доступных архивов для анализа практическую пользу имеет только архив sr.7z, весящий 120 Гб, после распаковки имеем 160 Гб данных (структура каталогов приведена на рисунке далее). С помощью rustyNog найти что-то интересное становится достаточно просто.



В логе каждая из записей содержит поля, позволяющие однозначно понять что было найдено и где (путь до файла, номер строки, выявленный тип данных). Список используемых по умолчанию полей доступен здесь. По умолчанию он содержит большое количество регулярных выражений, с помощью которых могут быть найдены секреты без применения энтропии.

Результат.

Пример выявленных интересных строк, с использованием энтропии приведён на рисунке ниже. Полностью результаты намеренно не привожу, т.к. анализ этих результатов показал отсутствие какой-либо полезной информации.

```

27175     {
27176         "stringsFound": [
27177             "R6flPri9fcebNaBhlzpBdRfMK5Z3KpIhHtmVdiBnaM8Nvd/WHwLqmuLMc3GkL30SgLDTMEZeS1SZ"
27178         ],
27179         "path": "./utility/BuildPatchTool/Engine/Content/Certificates/ThirdParty/cacert.pem",
27180         "reason": "Entropy",
27181         "linenum": 3274,
27182         "diff": "R6flPri9fcebNaBhlzpBdRfMK5Z3KpIhHtmVdiBnaM8Nvd/WHwLqmuLMc3GkL30SgLDTMEZeS1SZ\r"
27183     },
27184     {
27185         "stringsFound": [
27186             "lDePSHFjIuwXZQeVikvfj8ZaCuWw419eaxGrDPmF60Tp+ARz8un+XJiM9X0va7R+zdRcAitM0eGy"
27187         ],
27188         "path": "./utility/BuildPatchTool/Engine/Content/Certificates/ThirdParty/cacert.pem",
27189         "reason": "Entropy",
27190         "linenum": 352,
27191         "diff": "lDePSHFjIuwXZQeVikvfj8ZaCuWw419eaxGrDPmF60Tp+ARz8un+XJiM9X0va7R+zdRcAitM0eGy\r"
27192     },

```

Из полного списка шаблонов регулярных выражений для поиска заслуживают внимания следующие типы секретов: Generic Secret, RSA private key, Credentials in absolute URL, Generic Account API Key (остальные не рассматриваем, т.к. их наличие маловероятно). Результаты выявленных секретов приведены далее на рисунках.



```

155070     },
155071     {
155072         "stringsFound": [
155073             "Secret = \"5fe9feaf7ad2ee2080764faf8a5a45461eed231672c709cf8536aade8c0699eb\""
155074         ],
155075         "path": "./engine/config/platform/pc/serverconfig.ini",
155076         "reason": "Generic Secret",
155077         "linenum": 10,
155078         "diff": "ClientSecret = \"5fe9feaf7ad2ee2080764faf8a5a45461eed231672c709cf8536aade8c0699eb\""
155079     },

```

```

124569     {
124570         "stringsFound": [
124571             "-----BEGIN RSA PRIVATE KEY-----"
124572         ],
124573         "path": "./dev/external/libmicrohttpd/libmicrohttpd-0.9.64/src/testcurl/https/tls_test_keys.h",
124574         "reason": "RSA private key",
124575         "linenum": 26,
124576         "diff": "const char ca_key_pem[] = \"-----BEGIN RSA PRIVATE KEY-----\\n\\n\""
124577     },

```

```

89582     {
89583         "stringsFound": [
89584             "ftp://login:pass@bleh.gif"
89585         ],
89586         "path": "./dev/external/icu/source/test/testdata/regextst.txt",
89587         "reason": "Credentials in absolute URL",
89588         "linenum": 1818,
89589         "diff": "\"\\([Ii][Mm][Gg]\\)(\\S+?)\\(\\([Ii][Mm][Gg]\\)\" G \"<[img]ftp://login:pass@bleh.gif[/img]</>\"
89590     }

```

Итог анализа

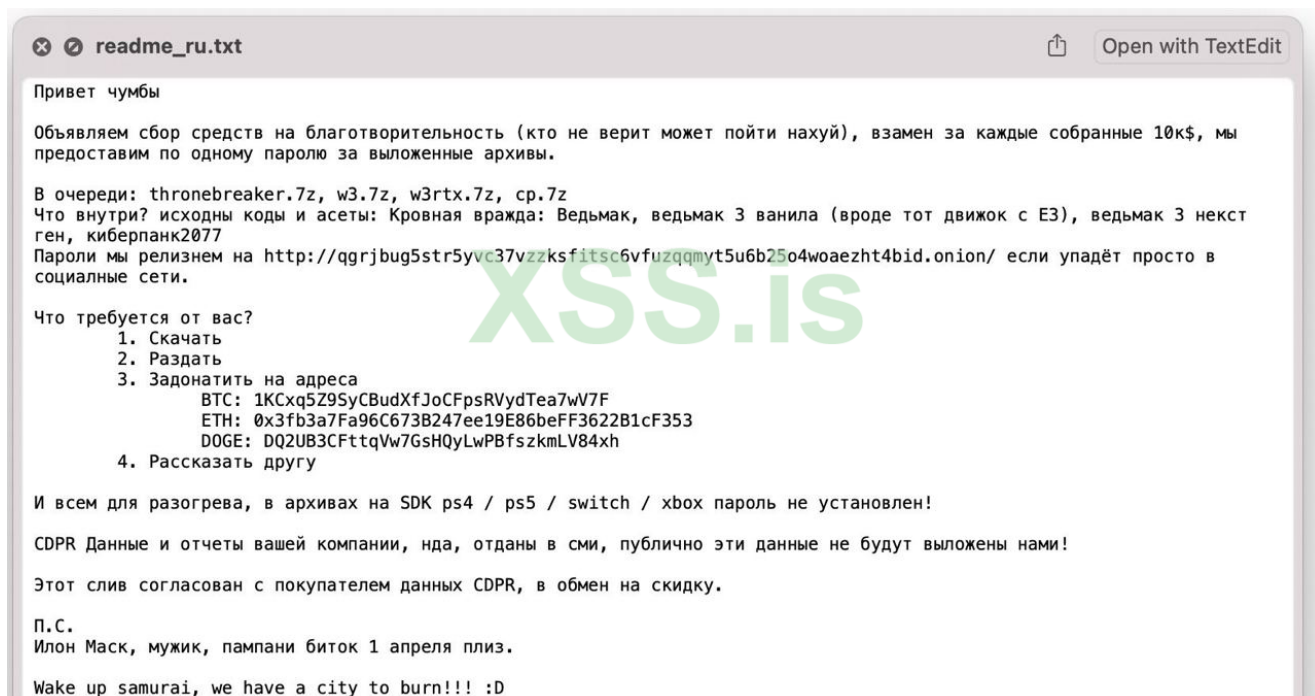
В составе "утечки" CD-Project слитой в паблик было выявлено:

- несколько частных ключей;
- несколько секретов, используемых при сборке-отладке;
- большое количество адресов email разработчиков используемых opensource компонентов;
- очень большая куча мусора.

Ничего из выявленного не представляет ценности как таковой.

Результат в целом соответствует действительности: были слиты части SDK и куски кода, не содержащие какой-либо ценности.

Всё то же, что имеет хоть какой-то вес содержится в архивах под паролем, о чём было написано в readme-файле в торрент раздаче данной "утечки".



Заключение

В этой статье мы рассмотрели простой способ для быстрого анализа данных, а вот где применять его решать вам.

Спасибо за внимание!

P.S.: Статя написана при содействии club1337, спасибо.

Last edited: Jul 15, 2021