



CyberDefense

MWS以外のデータセット紹介

MWS2023プレミーティング

2023/06/29

サイバーディフェンス研究所

中島将太

Malware関連のデータセット

- 有名なデータセット
 - Microsoft Malware Classification Challenge (BIG 2015)
 - EMBER dataset
 - SoReL-20M
- 興味深いデータセット
 - APT-Malware Dataset
 - Malpedia
- Google Dataset Search
- 研究利用のためのマルウェアの収集方法
 - VirusTotal
 - 無償利用可能サービス

有名なデータセット

Microsoft Malware Classification Challenge (BIG 2015)

- Microsoftが2015年にKaggleで開催したマルウェアの分類コンテストのデータセット
- 以下のマルウェアファミリーが含まれる
 1. Ramnit
 2. Lollipop
 3. Kelihos_ver3
 4. Vundo
 5. Simda
 6. Tracur
 7. Kelihos_ver1
 8. Obfuscator.ACY
 9. Gatak

The screenshot shows the Kaggle competition page for the Microsoft Malware Classification Challenge (BIG 2015). The page header includes the competition title, a prize money of \$16,000, and the host Microsoft. The main content area features a 'Description' section with a 'March 2018 Update' and a 'Timeline' section. The update text states: 'when using this dataset, please cite <http://arxiv.org/abs/1802.10135>'. Below the text is a grid of 12 blue icons representing various malware families and security concepts.

Research Prediction Competition

Microsoft Malware Classification Challenge (BIG 2015)

Classify malware into families based on file content and characteristics

Microsoft · 377 teams · 8 years ago

\$16,000 Prize Money

Overview Data Code Discussion Leaderboard Rules Team Submissions **Late Submission** ...

Overview

Description

March 2018 Update:
when using this dataset, please cite <http://arxiv.org/abs/1802.10135>

Evaluation

Prizes

Timeline

In recent years, the malware industry has become a well organized market involving large amounts of money. Well funded, multi-player syndicates invest heavily in technologies and capabilities built to evade traditional protection, requiring anti-malware vendors to develop counter mechanisms for finding and deactivating them. In the meantime, they inflict real financial and emotional pain to users of computer

Microsoft Malware Classification Challenge (BIG 2015)

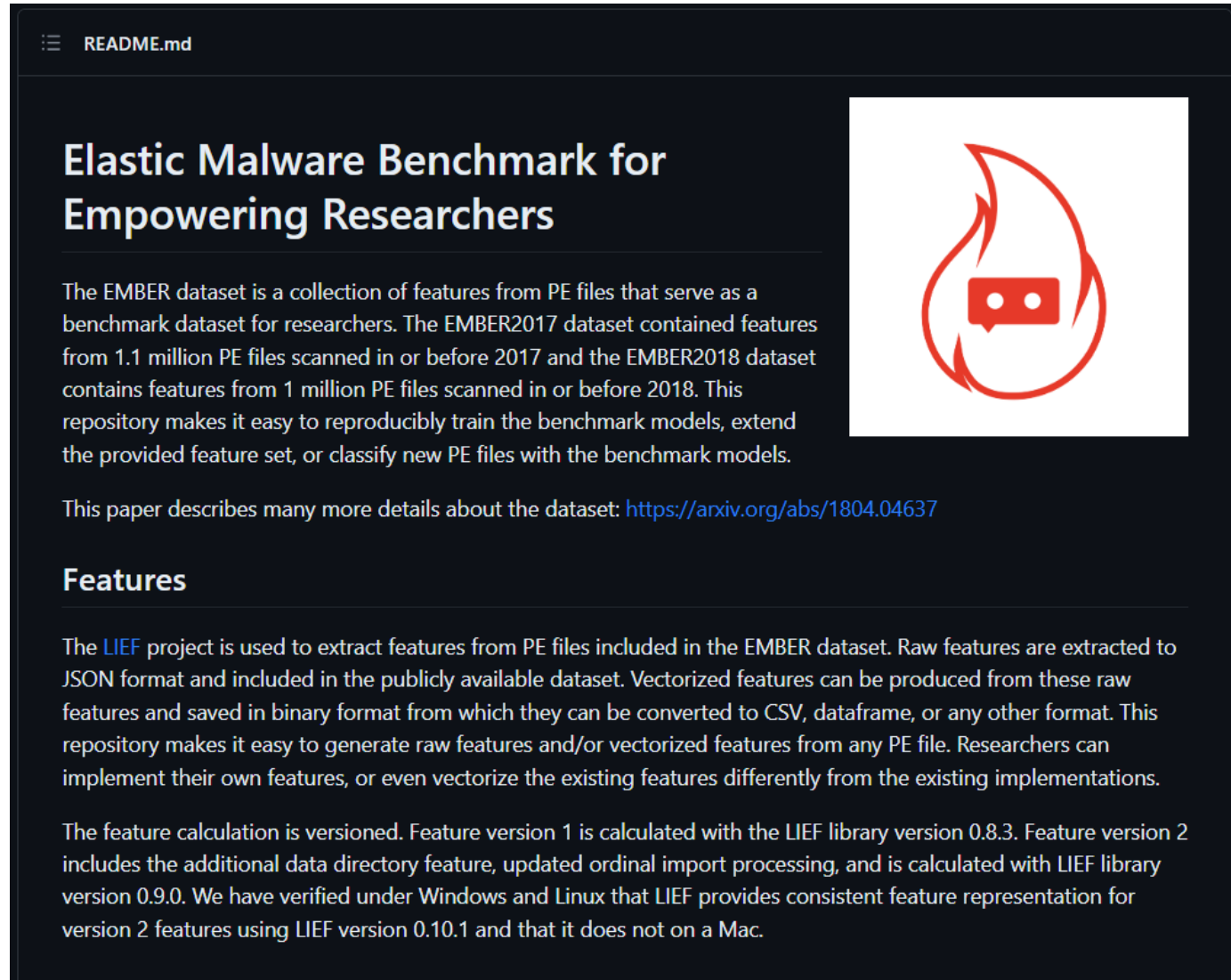
- PEヘッダを削ったHEXダンプとIDAの逆アセンブルデータが提供されている

```
1 00401000 56 8D 44 24 08 50 8B F1 E8 1C 1B 00 00 C7 06 08
2 00401010 BB 42 00 8B C6 5E C2 04 00 CC CC CC CC CC CC
3 00401020 C7 01 08 BB 42 00 E9 26 1C 00 00 CC CC CC CC CC
4 00401030 56 8B F1 C7 06 08 BB 42 00 E8 13 1C 00 00 F6 44
5 00401040 24 08 01 74 09 56 E8 6C 1E 00 00 83 C4 04 8B C6
6 00401050 5E C2 04 00 CC CC CC CC CC CC CC CC CC CC CC
7 00401060 8B 44 24 08 8A 08 8B 54 24 04 88 0A C3 CC CC CC
8 00401070 8B 44 24 04 8D 50 01 8A 08 40 84 C9 75 F9 2B C2
9 00401080 C3 CC CC CC CC CC CC CC CC CC CC CC CC CC CC
10 00401090 8B 44 24 10 8B 4C 24 0C 8B 54 24 08 56 8B 74 24
11 004010A0 08 50 51 52 56 E8 18 1E 00 00 83 C4 10 8B C6 5E
12 004010B0 C3 CC CC CC CC CC CC CC CC CC CC CC CC CC CC
13 004010C0 8B 44 24 10 8B 4C 24 0C 8B 54 24 08 56 8B 74 24
14 004010D0 08 50 51 52 56 E8 65 1E 00 00 83 C4 10 8B C6 5E
15 004010E0 C3 CC CC CC CC CC CC CC CC CC CC CC CC CC CC
16 004010F0 33 C0 C2 10 00 CC CC CC CC CC CC CC CC CC CC
17 00401100 B8 08 00 00 00 C2 04 00 CC CC CC CC CC CC CC
18 00401110 B8 03 00 00 00 C3 CC CC CC CC CC CC CC CC CC
19 00401120 B8 08 00 00 00 C3 CC CC CC CC CC CC CC CC CC
20 00401130 8B 44 24 04 A3 AC 49 52 00 B8 FE FF FF FF C2 04
21 00401140 00 CC CC CC CC CC CC CC CC CC CC CC CC CC CC
22 00401150 A1 AC 49 52 00 85 C0 74 16 8B 4C 24 08 8B 54 24
23 00401160 04 51 52 FF D0 C7 05 AC 49 52 00 00 00 00 00 B8
24 00401170 FB FF FF FF C2 08 00 CC CC CC CC CC CC CC CC
25 00401180 6A 04 68 00 10 00 00 68 68 BE 1C 00 6A 00 FF 15
26 00401190 9C 63 52 00 50 FF 15 C8 63 52 00 8B 4C 24 04 6A
27 004011A0 00 6A 40 68 68 BE 1C 00 50 89 01 FF 15 C4 63 52
28 004011B0 00 B8 04 00 00 00 C2 04 00 CC CC CC CC CC CC
```

```
03975 .text:00471DFA 89 75 E8          mov     [ebp+var_18], esi
03976 .text:00471DFD 5E          pop     esi
03977 .text:00471DFE BB 01 24 43 06      mov     ebx, 6432401h
03978 .text:00471E03 89 5D E8          mov     [ebp+var_18], ebx
03979 .text:00471E06 5B          pop     ebx
03980 .text:00471E07 8B 0D 0C 10 4B 00      mov     ecx, dword_4B100C
03981 .text:00471E0D 83 FA 87          cmp     edx, 0FFFFFF87h
03982 .text:00471E10 74 0B          jz      short loc_471E1D
03983 .text:00471E12 81 FA 60 18 7E 7E      cmp     edx, 7E7E1860h
03984 .text:00471E18 74 03          jz      short loc_471E1D
03985 .text:00471E1A 89 55 F8          mov     [ebp+var_8], edx
03986 .text:00471E1D
03987 .text:00471E1D          loc_471E1D:          ; CODE XREF: sub_471C4F+1C1Cmj
03988 .text:00471E1D          ; sub_471C4F+1C9Cmj
03989 .text:00471E1D 8B E5          mov     esp, ebp
03990 .text:00471E1F 5D          pop     ebp
03991 .text:00471E20 C2 04 00          retn   4
03992 .text:00471E20          sub_471C4F          endp
03993 .text:00471E20
03994 .text:00471E20          ; -----
03995 .text:00471E23 88          db 88h
03996 .text:00471E24 20 04 4C 08 21 12 08 04 40 02 00 00 88 01 44 41      dd 84C0420h, 4081221h, 240h, 41440188h
03997 .text:00471E34          ; [00000006 BYTES: COLLAPSED FUNCTION VirtualAlloc. PRESS KEYPAD CTRL- "+" TO EXPAND]
03998 .text:00471E3A          ; [00000006 BYTES: COLLAPSED FUNCTION GetCurrentThreadId. PRESS KEYPAD CTRL- "+" TO EXPAND]
03999 .text:00471E40          ; -----
04000 .text:00471E40 FF 25 1C F5 46 00      jmp     ds:_lcreat
04001 .text:00471E46          ; -----
04002 .text:00471E46 FF 25 E4 F4 46 00      jmp     ds:FindNextVolumeA
04003 .text:00471E4C          ; -----
04004 .text:00471E4C FF 25 E0 F4 46 00      jmp     ds:WriteFile
04005 .text:00471E52          ; -----
04006 .text:00471E52 FF 25 10 F5 46 00      jmp     ds:RegisterConsoleIME
04007 .text:00471E58          ; -----
04008 .text:00471E58 FF 25 24 F5 46 00      jmp     ds:GetBinaryTypeW
04009 .text:00471E5E          ; -----
04010 .text:00471E5E FF 25 18 F5 46 00      jmp     ds:BackupWrite
04011 .text:00471E64          ; -----
```

EMBER dataset

- Elasticが提供するデータセット
- LIEFで抽出した特徴量のみを提供している



☰ README.md

Elastic Malware Benchmark for Empowering Researchers


The EMBER dataset is a collection of features from PE files that serve as a benchmark dataset for researchers. The EMBER2017 dataset contained features from 1.1 million PE files scanned in or before 2017 and the EMBER2018 dataset contains features from 1 million PE files scanned in or before 2018. This repository makes it easy to reproducibly train the benchmark models, extend the provided feature set, or classify new PE files with the benchmark models.

This paper describes many more details about the dataset: <https://arxiv.org/abs/1804.04637>

Features

The [LIEF](#) project is used to extract features from PE files included in the EMBER dataset. Raw features are extracted to JSON format and included in the publicly available dataset. Vectorized features can be produced from these raw features and saved in binary format from which they can be converted to CSV, dataframe, or any other format. This repository makes it easy to generate raw features and/or vectorized features from any PE file. Researchers can implement their own features, or even vectorize the existing features differently from the existing implementations.

The feature calculation is versioned. Feature version 1 is calculated with the LIEF library version 0.8.3. Feature version 2 includes the additional data directory feature, updated ordinal import processing, and is calculated with LIEF library version 0.9.0. We have verified under Windows and Linux that LIEF provides consistent feature representation for version 2 features using LIEF version 0.10.1 and that it does not on a Mac.



- SophosとReversingLabsが公開したデータセット
- 200万のマルウェアサンプルが含まれる
 - データ量は8TB
- マルウェアの生バイナリ、学習済みのモデル、メタ情報を処理済みデータなどがまとめて公開されている

SoReL-20M

```
s3://sorel-20m/09-DEC-2020/
|   Terms and Conditions of Use.pdf -- the terms you agree to by using this data and code
|
+---baselines
|   +---checkpoints
|   |   +---FFNN - per-epoch checkpoints for 5 seeds of the feed-forward neural network
|   |   +---lightGBM - final trained lightGBM model for 5 seeds
|   |
|   +---results
|   |   ffnn_results.json - index file of results, required for plotting
|   |   lgbm_results.json - index file of results, required for plotting
|   |
|   +---FFNN
|   |   +---seed0-seed4 - individual seed results, ~1GB each
|   |
|   +---lightgbm
|   |   +---seed0-seed4 - individual seed results, ~1GB each
|
+---binaries
|   approximately 8TB of zlib compressed malware binaries
|
+---lightGBM-features
|   test-features.npz - array of test data for lightGBM; 37GB
|   train-features.npz - array of training data for lightGBM; 113GB
|   validation-features.npz - array of validation data for lightGBM; 22GB
|
+---processed-data
|   meta.db - contains index, labels, tags, and counts for the data; 3.5GB
|
+---ember_features - LMDB directory with baseline features, ~72GB
+---pe_metadata - LMDB directory with full metadata dumps, ~480GB
```

興味深いデータセット

APT Malware dataset

- マルウェアがAPTのマルウェアであるかを予測する分類器を作成する研究
- 19アクターの24ファミリー、3131サンプルを学習
- 論文で利用したデータセットを公開している

Table 1: APTs elements per classes

Class	Count	Class	Count
Patchwork	559	Lazarus Group	58
APT29	205	Sandwork	44
Winnti Group	176	Hurricane Panda	315
Carbanak	105	APT30	101
Volatile Cedar	35	Violin Panda	23
NSA	13	Desert Falcon	45
Transparent Tribe	267	Molerats	25
Shiqiang	31	APT28	68
Roaming Tiger	14	Lotus Blossom	48
Mirage	54		

APT-Malware

Introduced by Laurenza et al. in [Malware triage for early identification of Advanced Persistent Threat activities](#)

The APT Malware dataset is used to train classifiers to predict if a given malware belongs to the "Advanced Persistent Threat" (APT) type or not. It contains 3131 samples spread over 24 different unique malware classes.

Source: <https://arxiv.org/pdf/1810.07321.pdf>

[Homepage](#)

License ⓘ

Unknown

Modalities ⓘ

APT Malware dataset

- PEFrameを使ってマルウェアの静的な情報から4000以上の特徴を抽出して7クラスを作成

Optional Header (30 features). Every file has an optional header that provides information to the loader.

This header is optional in the sense that some files (specifically, object files) do not have it. For image files, this header is required. An object file can have an optional header, but generally this header has no function in an object file except to increase its size. Features are extracted from the optional header of the PE and contain information about the logical layout of the PE file, such as the address of the entry point, the alignment of sections, and the sizes of part of the file in memory.

MS-DOS Header (17 features). The MS-DOS executable-file header is composed of four distinct parts: a collection of header information (such as the signature word, the file size, etc.), a reserved section, a pointer to a Windows header (if one exists), and a stub program. MS-DOS uses the stub program to display a message if Windows has not been loaded when the user attempts to run a program. In this context, we are interested in features related to the execution of the file, including the number of bytes in the last page of the file, the number of pages, or the starting address of the Relocation Table.

File Header (18 features). The Windows executable-file header contains information that the loader requires for segmented executable files. This includes the linker version number, data specified by the linker, data specified by the resource compiler, tables of segment data, and tables of resource data. Moreover, the features related to this class highlight information about timestamp and the CPU platform that the PE is intended for.

Obfuscated String Statistics (3 features). Binaries contain program messages stored as strings that can be useful to understand their behavior. Classical tools like *String* extract byte sequences that can be readable strings to find these messages. Malware authors encode strings in their program to avoid extraction, in fact, even simple schemes can defeat this kind of tool and complicate static and dynamic analysis. In addition to PEFrame, we use functionalities of the FireEye Labs Obfuscated String Solver (*FLOSS*³). It is an open-source tool that automatically detects, extracts, and decodes obfuscated strings, such as malicious domains,

IP addresses, suspicious file paths, and so on, from Windows Portable Executable files availing of advanced static analysis techniques. We leverage this tool to compute some statistics, such as how many entry-points or relocations are present in the file.

Mutex (7 features). Mutex are objects commonly used to avoid simultaneous access to a resource, like a variable. If different software checks for the same mutex, then they can be linked. Our features are Boolean values that map the use of particular mutex identified in the training data.

Packer (64 features). Packers are software that compress binaries, keeping them executable. Similarly to the mutex related features, our features highlight if some particular packers are recognized.

Imported API (3917 features). Each software imports functions from common libraries or external files. The combination of imported functions can show similar behavior. We store this information as a vector with a Boolean value for each imported function, using the list of functions present in the training set as a taxonomy.

Buckets (98 features). Similar size in functions and directories can be a proof of similarity in the file structure and thus in the behavior. We observed that, usually, function size values range from 0 to 1,822 bytes, while directory size values range from 0 to 2,638 kbytes. To track these properties, both of them are subdivided into 49 buckets. Each bucket represents a range and counts the elements whose size is in the range. To choose the different ranges we observe the distribution of sizes and lengths in the training set, trying to form buckets that can better characterize the various classes.

APT Malware dataset

File Details

- *dataset.tar.gz* contains two hdf files containing features of APT-malware and normal malware
- *test_article.py* contains the code to test each implementation, included computation of some metrics
- *Checking_Result.py* contains the code to help computing quality and time metrics
- *ThresholdRandomForest.py* is a separate file containing all the methods to implement the functionalities of the first work
- *select** contain data obtained from our tests about best classes, best parameters and best columns to reproduce the published results

The screenshot displays a software interface with two main windows. The top window shows the file structure of 'malware_apt.h5', with 'class_df' expanded to show 'axis0', 'axis1', 'block0_items', 'block0_values', 'block1_items', 'block1_values', 'block2_items', and 'block2_values'. The 'block1_items' file is selected, and its 'Object Attribute Info' is shown on the right. The 'Attribute Creation Order' is 'Creation Order NOT Tracked', and there are 7 attributes. The attributes are listed in a table below.

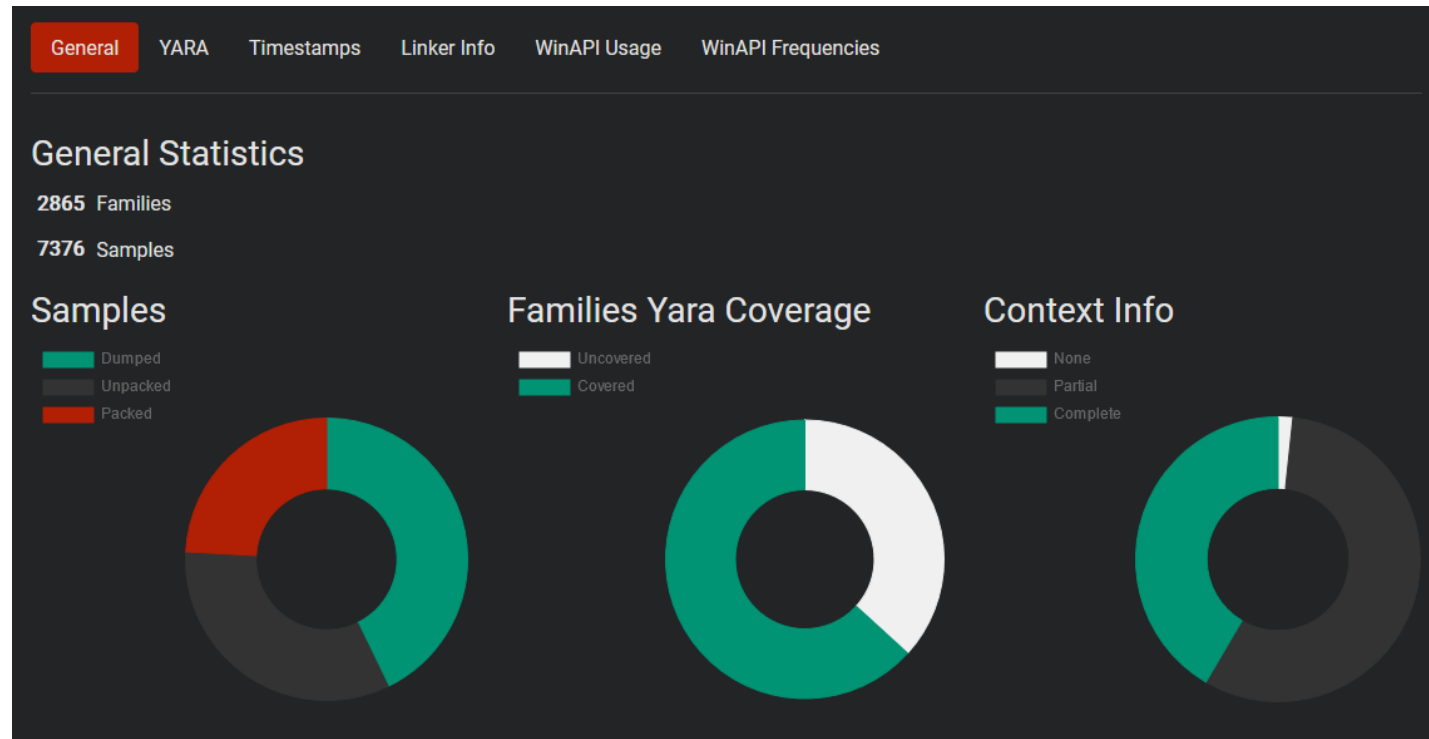
Name	Type	Array Size	Value[50](...)
CLASS	String, length = 5, padding = H5T_STR_NULLTERM, cset = H5T_CSET_UTF8	Scalar	ARRAY
FLAVOR	String, length = 5, padding = H5T_STR_NULLTERM, cset = H5T_CSET_UTF8	Scalar	numpy
TITLE	String, length = 1, padding = H5T_STR_NULLTERM, cset = H5T_CSET_UTF8	Scalar	
VERSION	String, length = 3, padding = H5T_STR_NULLTERM, cset = H5T_CSET_UTF8	Scalar	2.4
kind	String, length = 6, padding = H5T_STR_NULLTERM, cset = H5T_CSET_UTF8	Scalar	string
name	String, length = 2, padding = H5T_STR_NULLTERM, cset = H5T_CSET_ASCII	Scalar	N.
transposed	8-bit bitfield	Scalar	1

The bottom window, titled 'block1_items at /class_df/ [malware_apt.h5 in C:\Users\snakajima\Downloads\dataset.tar\dat...', shows a 'Data Display' window. It has a '0-based' index and a table with the following data:

Index	Value
0	sections_number
1	process32firstw
2	outputdebugstringa

Malpedia

- ドイツのFraunhofer FKIE (通信・情報処理・人間工学研究所) が運営するマルウェア情報集約サービス
 - ユーザが有益なブログや検体、YARAルールを投稿
 - ログインするとマルウェアのダウンロードも可能



Google Dataset Search

Google Dataset Search

Google



ログイン

Dataset Search

データセットを検索できます



[coronavirus covid-19](#) または [water quality site:canada.ca](#) をお試しください

データセット検索について [詳細をご確認ください](#)。

Google Dataset Search


The screenshot shows the Google Dataset Search interface. The search bar contains the word "malware". Below the search bar, there are filters for "最終更新日" (Last updated), "ダウンロード形式" (Download format), "ライセンス" (License), and "トピック" (Topic). There are also buttons for "すべて" (All), "1か月以内" (Within 1 month), "1年以内" (Within 1 year), and "過去3年間" (Past 3 years). A button for "保存済みのデータセット" (Saved datasets) is visible on the right.

The search results show 100+ datasets. The first result is "Windows PE Malware Dataset" by paperswithcode.com, updated on May 5, 2019. The second result is "Android Malware Dataset for Machine Learning" by kaggle.com, updated on Mar 13, 2021. The third result is "Windows Malware Detection Dataset" by figshare.com, updated on Mar 15, 2023. This result is highlighted with a red box. The fourth result is "Malware API Call Dataset" by ieee-dataport.org, updated on May 17, 2022. The fifth result is "Android malware dataset for machine learning 2" by figshare.com.

The detailed view of the "Windows Malware Detection Dataset" is shown on the right. It includes a "確認する" (Verify) button for figshare.com. A red box highlights the text "このデータセットを引用している学術記事: 327件 (Google Scholar で表示)". Below this, there is a "txt" file icon and a "一意の識別子" (Unique identifier) section with the URL <https://doi.org/10.6084/m9.figshare.21608262.v1>. Another red box highlights the "データセット更新日" (Dataset update date) section, which shows "Mar 15, 2023". Below this, there are sections for "データセットの提供元" (Dataset provider: figshare), "データセットの作成元" (Dataset creator: Irfan Yousuf), "ライセンス" (License: Attribution 4.0 (CC BY 4.0)), and "説明" (Description). The description states: "A dataset for Windows Portable Executable Samples with four feature sets. It contains four CSV files, one CSV file per feature set. 1. First feature set (DLLs_Imported.csv file) contains the DLLs imported by each malware family. The first column contains SHA256 values, second column contains the label or family type of the malware while the remaining columns list the names of imported DLLs. 2. Second feature set (API_Functions.csv files) contains the API functions called by these malware alongwith their SHA256 hash values and labels. 3. Third feature set (PE_Header.csv) contains values of 52 fields of PE header. All the fields are labelled in the CSV file. 4. Fourth feature set (PE_Section.csv file) contains 9 field values of 10 different PE sections. All the fields are labelled in the CSV file." Below the description, there is a "Malware Type / family Labels:" section with a legend: "0=Benign 1=RedLineStealer 2= Downloader 3=RAT 4=BankingTrojan 5=SnakeKeyLogger 6=Spyware".

研究利用のためのマルウェアの収集方法

VirusTotal

- 業界のデファクトスタンダードなサービス
- 有償契約ユーザーのみダウンロード可能 
 - APIの提供



弊社も販売パートナーです



VirusTotal

- 検知名を正解ラベルを利用可能
 - 正解ラベルとして利用すべきかの議論はまた別の話・・・
- 最近ベンダーの検知名を丸めたラベルを公式が提供

The screenshot shows the VirusTotal interface for a specific threat. At the top, it displays the 'Popular threat label' as 'trojan.tepfer/datastealer', 'Threat categories' as 'trojan', and 'Family labels' as 'tepfer', 'datastealer', and 'stealer'. Below this is a table titled 'Security vendors' analysis' with a sub-header 'Do you want to automate checks?'. The table lists various security vendors and their corresponding detection labels for the threat.

Security vendors' analysis		Do you want to automate checks?	
Acronis (Static ML)	⚠ Suspicious	AhnLab-V3	⚠ Trojan/Win32.Tepfer.R144050
Alibaba	⚠ TrojanPSW:Win32/Tepfer.966e38ff	ALYac	⚠ Generic.DataStealer.1.8F5D15EA
Antiy-AVL	⚠ Trojan[PSW]/Win32.Tepfer	Arcabit	⚠ Generic.DataStealer.1.8F5D15EA
Avast	⚠ Sf.Crypt-BI [Trj]	AVG	⚠ Sf.Crypt-BI [Trj]
Baidu	⚠ Win32.Trojan-PSW.Fareit.a	BitDefender	⚠ Generic.DataStealer.1.8F5D15EA
Bkav Pro	⚠ W32.AIDetectMalware	ClamAV	⚠ Win.Trojan.PonyStealer-9831667-0
CrowdStrike Falcon	⚠ Win/malicious_confidence_100% (W)	Cylance	⚠ Unsafe

無償利用可能な公開サービス

マルウェア共有サービス/IoC共有サービス

マルウェアのサンプルやIoCを共有するためのサービスがいくつかあります。もっとも登録が多く有名なサービスはVirusTotalですが、有償アカウントのみダウンロードが可能であるため、ここでは無償で利用可能なサービスを紹介します。

- [malpedia](#)
- [MalwareBazaar](#)
- [URLhaus](#)
- [ThreatFox](#)
- [MalShare](#)
- [VirusShare](#)
- [VirusBay](#)
- [Vx Vault](#)
- [theZoo](#)
- [vx-underground](#)

オンラインサンドボックス

無料で使用できるオンラインサンドボックスのサービスは本来マルウェアの挙動を解析するサービスですが、他者の投稿したマルウェアをダウンロードすることができます。

- [Hybrid Analysis](#)
- [Any Run](#)
- [Triage](#)
- [Joe Sandbox](#)
- [cape](#)
- [Cuckoo Sandbox](#)

参考URL

- <https://www.kaggle.com/competitions/malware-classification/overview>
- <https://github.com/elastic/ember>
- <https://paperswithcode.com/dataset/apt-malware>
- https://github.com/GiuseppeLaurenza/I_F_Identifier
- <https://malpedia.caad.fkie.fraunhofer.de/>
- https://github.com/pinksawtooth/how_to_become_a_malware_analyst

**ONLY
HUMANS CAN
COUNTER
HUMAN-DRIVEN
THREATS**