



LENNY ZELTSER



Analyzing Malicious Documents Cheat Sheet

General Approach to Document Analysis

1. Examine the document for anomalies, such as risky tags, scripts, or other anomalous aspects.
2. Locate embedded code, such as shellcode, VBA macros, JavaScript or other suspicious objects.
3. Extract suspicious code or object from the file.
4. If relevant, deobfuscate and examine JavaScript or macro code.
5. If relevant, disassemble and/or debug shellcode.
6. Understand the next steps in the infection chain.

Microsoft Office Format Notes

Binary document files supported by Microsoft Office use the OLE2 (a.k.a. Structured Storage) format.

SRP streams in OLE2 documents sometimes store a cached version of [earlier macro code](#).

```
pass to create file2.docm.
pcodedump.py -d file.doc Disassemble p-code macro code from file.doc.
rtfobj.py file.rtf Extract objects embedded into RTF-formatted file.rtf.
rtfdump.py file.rtf List groups and structure of RTF-formatted file.rtf.
rtfdump.py file.rtf -f 0 List groups in file.rtf that enclose an object.
rtfdump.py file.rtf -s 5 -H -d >out.bin Extract object from group 5 and save it into out.bin.
pyxsfw.py -xo file.doc Extract Flash (SWF) objects from OLE2 file file.doc.
```

Risky PDF Format Tags

/OpenAction and /AA specify the script or action to run automatically.

```
swf_mastah.py Extract Flash (SWF) objects from file.pdf into the out directory.
swf_mastah.py -f file.pdf -o out
```

Shellcode and Other Analysis Commands

```
xorsearch -W -d 3 file.bin Locate shellcode present in the binary file file.bin.
scdbg file.bin /foff 0x2B Emulate execution of shellcode in file.bin starting at offset 0x2B.
shellcode2exe file.bin Generate PE executable that runs shellcode.
base64dump.py file.txt List Base64-encoded data present in file file.txt.
base64dump.py file.txt -e bu Convert Base64-encoded data from file.txt as follows:
-s 2 -d >file.bin
```

This cheat sheet outlines tips and tools for analyzing malicious documents, such as Microsoft Office, RTF and Adobe Acrobat (PDF) files. To print it, use the one-page [PDF](#) version; you can also edit the [Word](#) version to customize it for your own needs.

General Approach to Document Analysis

1. Examine the document for anomalies, such as risky tags, scripts, or other anomalous

aspects.

2. Locate embedded code, such as shellcode, VBA macros, JavaScript or other suspicious objects.
3. Extract suspicious code or object from the file.
4. If relevant, deobfuscate and examine JavaScript or macro code.
5. If relevant, disassemble and/or debug shellcode.
6. Understand the next steps in the infection chain.

Microsoft Office Format Notes

- Binary document files supported by Microsoft Office use the OLE2 (a.k.a. Structured Storage) format.
- SRP streams in OLE2 documents sometimes store a cached version of [earlier macro code](#).
- OOXML documents (.docx, .xlsm, etc.) supported by MS Office use zip compression to store contents.
- Macros embedded in OOXML files are stored inside the OLE2 binary file, which is

within the zip archive.

- RTF documents don't support macros, but can contain other files embedded as OLE1 objects.

Useful MS Office File Analysis Commands

<code>unzip file.pptx</code>	Extract contents of OOXML file <i>file.pptx</i> .
<code>olevba.py file.xlsm</code> <code>olevba.py file.doc</code>	Locate and extract macros from <i>file.xlsm</i> or <i>file.doc</i> .
<code>oledump.py file.xls</code>	List all OLE2 streams present in <i>file.xls</i> .
<code>oledump.py -s 3 -v file.xls</code>	Extract macros stored inside stream 3 in <i>file.xls</i> .
<code>oledump.py file.xls -p plugin_http_heuristics</code>	Find obfuscated URLs in <i>file.xls</i> macros.
<code>msoffice-crypt -d -p pass file.docm</code> <code>file2.docm</code>	Decrypt OOXML file <i>file.docm</i> using password <i>pass</i> to create <i>file2.docm</i> .
<code>pcodedmp.py -d file.doc</code>	Disassemble p-code macro code from <i>file.doc</i> .

<code>rtfobj.py file.rtf</code>	Extract objects embedded into RTF-formatted <i>file.rtf</i> .
<code>rtfdump.py file.rtf</code>	List groups and structure of RTF-formatted <i>file.rtf</i> .
<code>rtfdump.py file.rtf -f O</code>	List groups in <i>file.rtf</i> that enclose an object.
<code>rtfdump.py file.rtf -s 5 -H -d > out.bin</code>	Extract object from group 5 and save it into <i>out.bin</i> .
<code>pyxswf.py -xo file.doc</code>	Extract Flash (SWF) objects from OLE2 file <i>file.doc</i> .

Risky PDF Format Tags

- `/OpenAction` and `/AA` specify the script or action to run automatically.
- `/JavaScript` and `/JS` specify JavaScript to run.
- `/GoTo` changes the view to a specified destination within the PDF or in another PDF file.
- `/Launch` can launch a program or open a document.

- /URI accesses a resource by its URL.
- /SubmitForm and /GoToR can send data to URL.
- /RichMedia can be used to embed Flash in a PDF.
- /ObjStm can hide objects inside an Object Stream.
- Be mindful of obfuscation with hex codes, such as /JavaScript vs. /J#61vaScript. ([See examples.](#))

Useful PDF File Analysis Commands

pdfid.py <i>file.pdf</i>	Scan <i>file.pdf</i> for risky keywords and dictionary entries.
peepdf.py -fl <i>file.pdf</i>	Examine <i>file.pdf</i> for risky tags and malformed objects.
pdf-parser.py – object <i>id file.pdf</i>	Display contents of object <i>id</i> in <i>file.pdf</i> . Add “–filter –raw” to decode the object’s stream.
qpdf – password= <i>pass</i> –decrypt <i>infile.pdf</i> <i>outfile.pdf</i>	Decrypt <i>infile.pdf</i> using password <i>pass</i> to create <i>outfile.pdf</i> .

`swf_mastah.py -f file.pdf -o out`

Extract Flash (SWF) objects from *file.pdf* into the *out* directory.

Shellcode and Other Analysis Commands

`xorsearch -W -d 3 file.bin`

Locate shellcode patterns inside the binary file *file.bin*.

`scdbg file.bin /foff 0x2B`

Emulate execution of shellcode in *file.bin* starting at offset *0x2B*.

`shellcode2exe file.bin`

Generate PE executable *file.exe* that runs shellcode from *file.bin*.

`jmp2it file.bin 0x2B`

Execute shellcode in file *file.bin* starting at offset *0x2B*.

`base64dump.py file.txt`

List Base64-encoded strings present in file *file.txt*.

`base64dump.py file.txt -e bu -s 2 -d > file.bin`

Convert backslash Unicode-encoded Base64 string #2 from *file.txt* as *file.bin* file.

Additional Document Analysis Tools

- [SpiderMonkey](#), [V8](#) and [box-js](#) help deobfuscate JavaScript that you extract from document files.
- [PDF Stream Dumper](#) combines several PDF analysis utilities under a single graphical user interface.
- [ViperMonkey](#) emulates VBA macro execution.
- [VirusTotal](#) and some [automated analysis sandboxes](#) can analyze aspects of malicious document files.
- [Hachoir-urwid](#) can display OLE2 stream contents.
- [101 Editor](#) (commercial) and [FileInsight](#) hex editors can parse and edit OLE structures.
- [ExeFilter](#) can filter scripts from Office and PDF files.
- [REMnux](#) distro includes many of the free document analysis tools mentioned above.

Post-Scriptum

Special thanks for feedback to [Pedro Bueno](#) and [Didier Stevens](#). If you have suggestions for improving this cheat sheet, please

let me know. [Creative Commons v3 “Attribution” License](#) for this cheat sheet version 3.0.

Updated September 1, 2017

MORE ON

[Information Security](#)

[Malicious Software](#)

The SANS [malware analysis course](#) I’ve co-authored explains the techniques summarized in this cheat sheet and covers many other reverse-engineering topics.

If you like this reference, take a look at my other [IT and security cheat sheets](#).

SHARE 

DID YOU LIKE THIS?

Follow me for more of the good stuff.



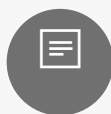
TWITTER



GOOGLE+



RSS FEED



NEWSLETTER

About the Author

Lenny Zeltser is a seasoned business and technology leader with extensive information security experience. He builds innovative endpoint defense solutions as VP of Products at [Minerva](#). He also trains incident response and digital forensics professionals at [SANS Institute](#). Lenny frequently speaks at industry events, writes articles and has co-authored books. He has earned the prestigious GIAC Security Expert designation, has an MBA from MIT Sloan and a Computer Science degree from the University of Pennsylvania.

[Learn more](#)