

# 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.

<code>pcodecmp.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 0</code>	List groups in <i>file.rtf</i> that enclose an object.
<code>rtfdump.py file.rtf -s 5 -H -d &gt;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.

<code>swf_mastah.py -f file.pdf -o out</code>	Extract Flash (SWF) objects from <i>file.pdf</i> into the output file <i>out</i> .
---	--

## Shellcode and Other Analysis Commands

<code>xorsearch -W -d 3 file.bin</code>	Locate shellcode patterns in the binary file <i>file.bin</i> .
<code>scdbg file.bin /foff 0x2B</code>	Emulate execution in <i>file.bin</i> starting at offset 0x2B.
<code>shellcode2exe file.bin</code>	Generate PE executable that runs shellcode from <i>file.bin</i> .
<code>base64dump.py file.txt</code>	List Base64-encoded strings present in file <i>file.txt</i> .
<code>base64dump.py file.txt -e bu -s 2 -d &gt;file.bin</code>	Convert backslash-encoded Base64 strings from <i>file.txt</i> as file <i>file.bin</i> .

## MORE ON

[Information Security](https://zeltser.com/information-security)

(<https://zeltser.com/information-security>)

[Malicious Software](https://zeltser.com/malicious-software)

(<https://zeltser.com/malicious-software>)

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 \(/media/docs/analyzing-malicious-document-files.pdf\)](/media/docs/analyzing-malicious-document-files.pdf) version; you can also edit the [Word \(/media/docs/analyzing-malicious-document-files.docx\)](/media/docs/analyzing-malicious-document-files.docx) 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.

The [SANS malware analysis course \(https://sans.org/for610\)](https://sans.org/for610) 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 \(/cheat-sheets/\)](/cheat-sheets/).

[SHARE](#) ➞

## 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 (<https://digital-forensics.sans.org/blog/2014/06/05/srp-streams-in-office-documents-reveal-earlier-macros>).
- 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

unzip <i>file.pptx</i>	Extract contents of OOXML file <i>file.pptx</i> .
olevba.py ( <a href="https://github.com/decalage2/oletools/wiki/olevba">https://github.com/decalage2/oletools/wiki/olevba</a> ) <i>file.xlsm</i> olevba.py <i>file.doc</i>	Locate and extract macros from <i>file.xlsm</i> or <i>file.doc</i> .
oledump.py ( <a href="https://blog.didierstevens.com/programs/oledump-py/">https://blog.didierstevens.com/programs/oledump-py/</a> ) <i>file.xls</i>	List all OLE2 streams present in <i>file.xls</i> .
oledump.py -s 3 -v <i>file.xls</i>	Extract macros stored inside stream 3 in <i>file.xls</i> .

<code>rtfdump.py file.xls -p plugin_http_heuristics</code> <code>LENNY ZELTSEY (HTTPS://ZELTSEY.COM/)</code>	Find obfuscated URLs in <i>file.xls</i> macros.
<code>msoffice-crypt (https://github.com/herumi/msoffice)</code> <code>-d -p pass file.docm file2.docm</code>	Decrypt OOOXML file <i>file.docm</i> using password <i>pass</i> to create <i>file2.docm</i> .
<code>pcodedmp.py</code> <code>(https://github.com/bontchev/pcodedmp) -d file.doc</code>	Disassemble p-code macro code from <i>file.doc</i> .
<code>rtfobj.py (https://www.decorage.info/python/rtfobj)</code> <code>file.rtf</code>	Extract objects embedded into RTF-formatted <i>file.rtf</i> .
<code>rtfdump.py</code> <code>(https://blog.didierstevens.com/2016/08/02/rtfdump-update-and-videos/) 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.

`#!python file.py s 5 H -o out.bin`  
LENNY ZELT SER (HTTPS://ZELTSER.COM/)

Extract  
object from  
group 5 and  
save it into  
*out.bin*.

`pyxswf.py`  
(<https://www.decalage.info/python/pyxswf>) -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.
- /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 (<https://blog.didierstevens.com/2008/04/29/pdf-let-me-count-the-ways/>)).)

## Useful PDF File Analysis Commands

`pdfid.py`  
(<https://blog.didierstevens.com/programs/pdf-tools/>) *file.pdf*  
Scan *file.pdf* for  
risky keywords and  
dictionary entries.

`peepdf.py` (<http://eternal-todo.com/tools/peepdf-pdf-analysis-tool/>) -fl  
*file.pdf*  
Examine *file.pdf* for  
risky tags and  
malformed objects.

<a href="https://zeltser.com/">pdf-parser.py</a> (ZELTSER (HTTPS://ZELTSER.COM/)) <a href="https://blog.didierstevens.com/programs/pdf-tools/">(https://blog.didierstevens.com/programs/pdf-tools/)</a> --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.
<a href="http://qpdf.sourceforge.net/">qpdf</a> ( <a href="http://qpdf.sourceforge.net/">http://qpdf.sourceforge.net/</a> ) --password= <i>pass</i> --decrypt <i>infile.pdf outfile.pdf</i>	Decrypt <i>infile.pdf</i> using password <i>pass</i> to create <i>outfile.pdf</i> .
<a href="https://zeltser.com/extracting-swf-from-pdf-using-swf-mastah/">swf_mastah.py</a> ( <a href="https://zeltser.com/extracting-swf-from-pdf-using-swf-mastah/">https://zeltser.com/extracting-swf-from-pdf-using-swf-mastah/</a> ), -f <i>file.pdf</i> -o <i>out</i>	Extract Flash (SWF) objects from <i>file.pdf</i> into the <i>out</i> directory.

## Shellcode and Other Analysis Commands

<a href="https://blog.didierstevens.com/2014/09/29/update-xorsearch-with-shellcode-detector/">xorsearch</a> <a href="https://blog.didierstevens.com/2014/09/29/update-xorsearch-with-shellcode-detector/">(https://blog.didierstevens.com/2014/09/29/update-xorsearch-with-shellcode-detector/)</a> , -W -d 3 <i>file.bin</i>	Locate shellcode patterns inside the binary file <i>file.bin</i> .
<a href="http://sandsprite.com/blogs/index.php?uid=7&amp;pid=152">scdbg</a> ( <a href="http://sandsprite.com/blogs/index.php?uid=7&amp;pid=152">http://sandsprite.com/blogs/index.php?uid=7&amp;pid=152</a> ) <i>file.bin</i> /foff 0x2B	Emulate execution of shellcode in <i>file.bin</i> starting at offset 0x2B.
<a href="https://zeltser.com/convert-shellcode-to-assembly/">shellcode2exe</a> ( <a href="https://zeltser.com/convert-shellcode-to-assembly/">https://zeltser.com/convert-shellcode-to-assembly/</a> ) <i>file.bin</i>	Generate PE executable <i>file.exe</i> that runs shellcode from <i>file.bin</i> .

<a href="https://digital-forensics.sans.org/blog/2014/12/30/taking-control-of-the-instruction-pointer/">https://digital-forensics.sans.org/blog/2014/12/30/taking-control-of-the-instruction-pointer/</a> <i>file.bin 0x2B</i>	Execute shellcode in file <i>file.bin</i> starting at offset <i>0x2B</i> .
<a href="https://blog.didierstevens.com/2017/07/02/update-base64dump-py-version-0-0-7/">https://blog.didierstevens.com/2017/07/02/update-base64dump-py-version-0-0-7/</a> <i>file.txt</i>	List Base64-encoded strings present in file <i>file.txt</i> .
<i>base64dump.py file.txt -e bu -s 2 -d &gt; file.bin</i>	Convert backslash Unicode-encoded Base64 string <i>#2</i> from <i>file.txt</i> as <i>file.bin</i> file.

## Additional Document Analysis Tools

- [SpiderMonkey](https://developer.mozilla.org/en-US/docs/Mozilla/Projects/SpiderMonkey) (<https://developer.mozilla.org/en-US/docs/Mozilla/Projects/SpiderMonkey>), [V8](https://isc.sans.edu/diary/V8+as+an+Alternative+to+SpiderMonkey+for+JavaScript+Deobfuscation/12157) (<https://isc.sans.edu/diary/V8+as+an+Alternative+to+SpiderMonkey+for+JavaScript+Deobfuscation/12157>) and [box-js](https://github.com/CapacitorSet/box-js) (<https://github.com/CapacitorSet/box-js>) help deobfuscate JavaScript that you extract from document files.
- [PDF Stream Dumper](https://zeltser.com/pdf-stream-dumper-malicious-file-analysis/) (<https://zeltser.com/pdf-stream-dumper-malicious-file-analysis/>) combines several PDF analysis utilities under a single graphical user interface.
- [ViperMonkey](https://github.com/decalage2/ViperMonkey) (<https://github.com/decalage2/ViperMonkey>) emulates VBA macro execution.
- [VirusTotal](https://www.virustotal.com/) (<https://www.virustotal.com/>) and some [automated analysis sandboxes](https://automated-malware-analysis.com/) ([/automated-malware-analysis/](https://automated-malware-analysis.com/)) can analyze aspects of malicious document files.
- [Hachoir-urwid](https://bitbucket.org/haypo/hachoir/wiki/hachoir-urwid) (<https://bitbucket.org/haypo/hachoir/wiki/hachoir-urwid>) can display OLE2 stream contents.

- **101 Editor** (<https://www.sweetscape.com/010editor/>) (commercial) and **LENNY ZELT SER** ([HTTPS://ZELT SER.COM/](https://zeltser.com/)) **FileInsight** (<https://www.mcafee.com/us/downloads/free-tools/fileinsight.aspx>) hex editors can parse and edit OLE structures.
- **ExeFilter** (<http://www.decalage.info/exefilter>) can filter scripts from Office and PDF files.
- **REMnux** (<https://remnux.org/>) distro includes many of the free document analysis tools mentioned above.



## Post-Scriptum

Special thanks for feedback to [Pedro Bueno](http://handlers.dshield.org/pbueno/) (<http://handlers.dshield.org/pbueno/>) and [Didier Stevens](http://blog.didierstevens.com/) (<http://blog.didierstevens.com/>). If you have suggestions for improving this cheat sheet, please [let me know \(/contact/\)](#). [Creative Commons v3 "Attribution" License](#) (<http://creativecommons.org/licenses/by/3.0/>) for this cheat sheet version 3.0.

*Updated September 6, 2017*

---

### DID YOU LIKE THIS?

Follow me for more of the good stuff.

---

## About the Author

Lenny Zeltser develops teams, products, and programs that use information security to achieve business results. He is presently the CISO at Axonius and an author and instructor at SANS Institute. Over the past two decades, Lenny has been leading efforts to establish resilient security practices and solve hard security problems. As a respected author and speaker, he has been advancing cybersecurity tradecraft and contributing to the community. His insights build upon 20 years of real-world experiences, a Computer Science degree from the University of Pennsylvania, and an MBA degree from MIT Sloan.

[Learn more \(https://zeltser.com/about\)](https://zeltser.com/about)

Copyright © 1995-2020 Lenny Zeltser. All rights reserved.