



ANALYZING MALICIOUS DOCUMENTS

This cheat sheet outlines tips and tools for analyzing malicious documents, such as Microsoft Office, RTF and Adobe Acrobat (PDF) files.

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>	Locate and extract macros from <i>file.xlsm</i> or <i>file.doc</i> .
<code>olevba.py file.doc</code>	

<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 file2.docm</code>	Decrypt OOXML file <i>file.docm</i> using password <i>pass</i> to create <i>file2.docm</i> .
<code>p-codedmp.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 >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

<code>pdfid.py file.pdf</code>	Scan <i>file.pdf</i> for risky keywords and dictionary entries.
<code>peepdf.py -f1 file.pdf</code>	Examine <i>file.pdf</i> for risky tags and malformed objects.

<code>pdf-parser.py --object id file.pdf</code>	Display contents of object <i>id</i> in <i>file.pdf</i> . Add " <code>--filter --raw</code> " to decode the object's stream.
<code>qpdf --password=pass --decrypt infile.pdf outfile.pdf</code>	Decrypt <i>infile.pdf</i> using password <i>pass</i> to create <i>outfile.pdf</i> .
<code>swf_mastah.py -f file.pdf -o out</code>	Extract Flash (SWF) objects from <i>file.pdf</i> into the <i>out</i> directory.

Shellcode and Other Analysis Commands

<code>xorsearch -W -d 3 file.bin</code>	Locate shellcode patterns inside the binary file <i>file.bin</i> .
<code>sctdbg file.bin /foff 0x2B</code>	Emulate execution of shellcode in <i>file.bin</i> starting at offset <code>0x2B</code> .
<code>shellcode2exe file.bin</code>	Generate PE executable <i>file.exe</i> that runs shellcode from <i>file.bin</i> .
<code>jmp2it file.bin 0x2B</code>	Execute shellcode in file <i>file.bin</i> starting at offset <code>0x2B</code> .
<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 >file.bin</code>	Convert backslash Unicode-encoded Base64 string <code>#2</code> from <i>file.txt</i> as <i>file.bin</i> 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.