

## 104 – PDF Classification

### Team Information

Team Name : LuckyVicky

Team Member : Eungchang Lee, Hyun Yi, Juho Heo, Dongkyu Lee

Email Address : dfc\_luckyvicky@googlegroups.com

### Instructions

**Description** Develop a tool to analyze the internal structure of given example PDF files and classify them into normal PDFs, corrupted PDFs, and encrypted PDFs.

Target	Hash (MD5)
1.pdf	c578c2eaeb340372f7ca96ae20c5d6e5
2.pdf	89a5c129333d6aba39befe9b6ce40cf7
3.pdf	d07f21d1a4ef17de5ec7a531f50b9f84

Teams must:

- Develop and document the step-by-step approach used to solve this problem to allow another examiner to replicate team actions and results.
- Specify all tools used in deriving the conclusion(s).

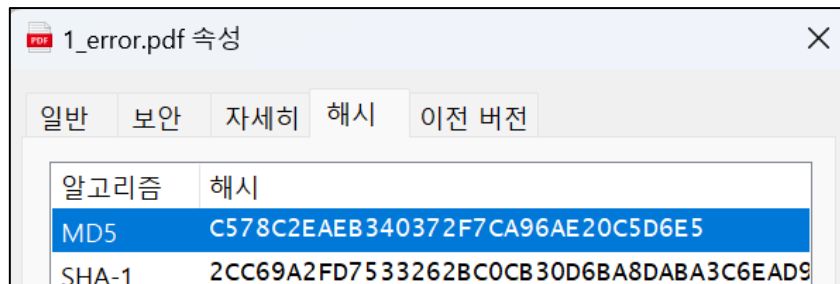
**Tools used:**

Name:	Python	Publisher:	Python Software Foundation
Version:	3.11.8		
URL:	<a href="https://www.python.org/">https://www.python.org/</a>		

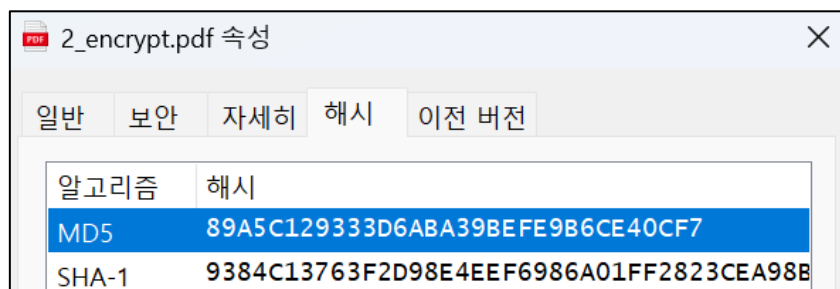
Name:	pypdf	Publisher:	py-pdf
Version:	4.3.1		
URL:	<a href="https://github.com/py-pdf/pypdf">https://github.com/py-pdf/pypdf</a>		

Name:	HxD	Publisher:	Maël Hörz
Version:	1.7.7.0		
URL:	<a href="https://mh-nexus.de/en/hxd/">https://mh-nexus.de/en/hxd/</a>		

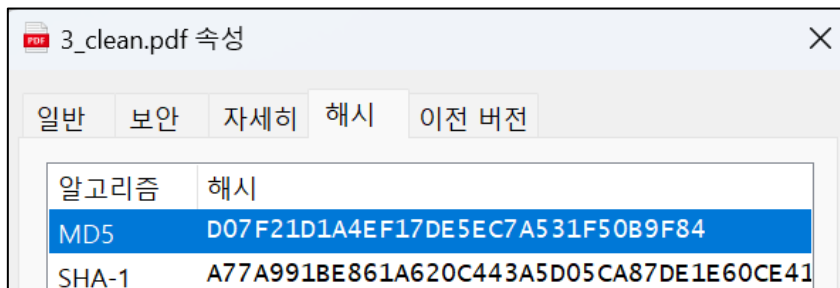
## Step-by-step methodology:



[그림 1] 1\_error.pdf 파일의 md5 해시 값 확인



[그림 2] 2\_encrypt.pdf 파일의 md5 해시 값 확인



[그림 3] 3\_clean.pdf 파일의 md5 해시 값 확인

문제로 주어진 파일들의 MD5 해시 값이 일치함을 확인하였습니다.

```

E:\dfc2024\104 - PDF Classification>dir /S
E 드라이브의 볼륨에는 이름이 없습니다.
볼륨 일련 번호: 3239-0852

E:\dfc2024\104 - PDF Classification 디렉터리

2024-08-16 오후 02:44 <DIR> .
2024-08-11 오후 08:52 <DIR> ..
2024-08-16 오후 02:44 72,213 104 - PDF Classification.docx
2024-08-16 오후 02:44 <DIR> pdfs
2024-08-16 오후 02:46 2,544 pdf_analyzer.py
                2개 파일                74,757 바이트

E:\dfc2024\104 - PDF Classification\pdfs 디렉터리

2024-08-16 오후 02:44 <DIR> .
2024-08-16 오후 02:44 <DIR> ..
2024-08-02 오전 02:50 704,228 1_error.pdf
2024-07-22 오후 04:37 687,547 2_encrypt.pdf
2024-07-22 오후 03:55 704,228 3_clean.pdf
                3개 파일                2,096,003 바이트

전체 파일:
        5개 파일                2,170,760 바이트
        5개 디렉터리    33,530,843,136 바이트 남음

```

[그림 4] 현재 폴더에 존재하는 파일 목록 출력

현재 폴더에는 pdf\_analyzer.py 파일과 pdfs 폴더에 3개의 pdf 파일이 존재합니다.

```

E:\dfc2024\104 - PDF Classification>python pdf_analyzer.py
usage: pdf_analyzer.py [-h] --path PATH
pdf_analyzer.py: error: the following arguments are required: --path

E:\dfc2024\104 - PDF Classification>python pdf_analyzer.py -h
usage: pdf_analyzer.py [-h] --path PATH

Analyze PDF files

options:
  -h, --help  show this help message and exit
  --path PATH  Path to the PDF file

```

[그림 5] pdf\_analyzer.py 실행 결과

pdf\_analyzer.py를 실행했을 때 필요한 argument를 알려주며, -h 옵션을 사용하면 상세한 옵션 설명을 출력해줍니다.

```
E:\dfc2024\104 - PDF Classification>python pdf_analyzer.py --path pdfs
File: 1_error.pdf
- Corrupted: True
- Encrypted: False
Total number of pages: 5

Page 1 Structure:
{'/Type': '/Page', '/Parent': IndirectObject(2, 0, 1786596359120), '/Resources': {'/ExtGState': {'/GS5': IndirectObject(5, 0, 1786596359120), '/GS8': IndirectObject(8, 0, 1786596359120)}, '/Font': {'/F1': IndirectObject(6, 0, 1786596359120), '/F2': IndirectObject(9, 0, 1786596359120)}, '/ProcSet': ['/PDF', '/Text', '/ImageB', '/ImageC', '/ImageI']}, '/MediaBox': [0, 0, 960, 540], '/Contents': IndirectObject(4, 0, 1786596359120), '/Group': {'/Type': '/Group', '/S': '/Transparency', '/CS': '/DeviceRGB'}, '/Tabs': '/S', '/StructParents': 0}
Key: /Type, Value: /Page
Key: /Parent, Value: {'/Type': '/Pages', '/Count': 8, '/Kids': [IndirectObject(3, 0, 1786596359120), IndirectObject(11, 0, 1786596359120), IndirectObject(20, 0, 1786596359120), IndirectObject(30, 0, 1786596359120), IndirectObject(38, 0, 1786596359120), IndirectObject(52, 0, 1786596359120), IndirectObject(60, 0, 1786596359120), IndirectObject(71, 0, 1786596359120)]}
Key: /Resources, Value: {'/ExtGState': {'/GS5': IndirectObject(5, 0, 1786596359120), '/GS8': IndirectObject(8, 0, 1786596359120)}, '/Font': {'/F1': IndirectObject(6, 0, 1786596359120), '/F2': IndirectObject(9, 0, 1786596359120)}, '/ProcSet': ['/PDF', '/Text', '/ImageB', '/ImageC', '/ImageI']}
Key: /MediaBox, Value: [0, 0, 960, 540]
Key: /Contents, Value: {'/Filter': '/FlateDecode'}
Key: /Group, Value: {'/Type': '/Group', '/S': '/Transparency', '/CS': '/DeviceRGB'}
Key: /Tabs, Value: /S
Key: /StructParents, Value: 0
```

#### [그림 6] pdf\_analyzer.py 분석 결과 - 1

--path 옵션에 pdf 파일들이 존재하는 경로를 지정한 후 코드를 실행하면 각 파일들의 손상 및 암호 설정 유무와 페이지 별 구조를 출력해줍니다.

```
File: 2_encrypt.pdf
- Corrupted: False
- Encrypted: True
The file is encrypted and cannot be analyzed.
```

#### [그림 7] pdf\_analyzer.py 분석 결과 - 2

암호화된 파일은 구조를 분석할 수 없기 때문에 위와 같은 문구를 출력합니다.

```
--- Summary Results ---
File: 1_error.pdf, Corrupted: True, Encrypted: False, Type: Corrupted PDF
File: 2_encrypt.pdf, Corrupted: False, Encrypted: True, Type: Encrypted PDF
File: 3_clean.pdf, Corrupted: False, Encrypted: False, Type: Normal PDF
```

#### [그림 8] pdf\_analyzer.py 분석 결과 - 3

마지막으로 분석 결과를 요약하여 문제에서 요구하는 분류 결과를 한 눈에 확인할 수 있으며, Python에서 실행할 수 있는 스크립트 파일과 라이브러리 설치를 위한 requirements.txt 파일을 압축하여 제출합니다.



```
# remove wrong objects (not pointing to correct structures) - cf #2326
if not self.strict:
    loc = stream.tell()
    for gen, xref_entry in self.xref.items():
        if gen == 65535:
            continue
        ids = list(xref_entry.keys())
        for id in ids:
            stream.seek(xref_entry[id], 0)
            try:
                self.read_object_header(stream)
            except ValueError:
                logger_warning(
                    f"Ignoring wrong pointing object {id} {gen} (offset {xref_entry[id]}",
                    __name__,
                )
                del xref_entry[id] # we can delete the id, we are parsing ids
    stream.seek(loc, 0) # return to where it was
```

[그림 12] pypdf에서 손상된 파일을 처리하는 코드<sup>2</sup>

[그림 12]의 코드를 보면 xref entry 별 header를 읽어오는 과정에서 에러가 발생할 경우 [그림 11]과 같은 에러가 발생하게 됩니다. [그림 12] 코드의 주석에 적혀 있듯이 object가 삭제되어 구조가 올바르지 않은 상태입니다.

00019460	BE 01 10 C7 6E D0 0D 0A 65 6E 64 73 74 72 65 61	%..ÇnD..endstrea
00019470	6D 0D 0A 65 6E 64 6F 62 6A 0D 0A 34 35 20 30 20	m..endobj...45 0
00019480	6F 62 6A 0D 0A 3C 3C 2F 54 79 70 65 2F 58 4F 62	obj...<</Type/XOb
00019490	6A 65 63 74 2F 53 75 62 74 79 70 65 2F 49 6D 61	ject/Subtype/Ima
000194A0	67 65 2F 57 69 64 74 68 20 35 34 31 2F 48 65 69	ge/Width 541/Hei
000194B0	67 68 74 20 37 37 34 2F 43 6F 6C 6F 72 53 70 61	ght 774/ColorSpa
000194C0	63 65 2F 44 65 76 69 63 65 52 47 42 2F 42 69 74	ce/DeviceRGB/Bit
000194D0	73 50 65 72 43 6F 6D 70 6F 6E 65 6E 74 20 38 2F	sPerComponent 8/
000194E0	49 6E 74 65 72 70 6F 6C 61 74 65 20 66 61 6C 73	Interpolate fals
000194F0	65 2F 46 69 6C 74 65 72 2F 46 6C 61 74 65 44 65	e/Filter/FlateDe
00019500	63 6F 64 65 2F 4C 65 6E 67 74 68 20 31 36 37 35	code/Length 1675
00019510	38 39 3E 3E 0D 0A 73 74 72 65 61 6D 0D 0A 78 9C	89>>..stream..xœ
00019520	EC 9D 09 7C 13 D5 F6 C7 23 50 76 4A 91 4D 28 4B	i.. .ÕöÇ#PvJ'M(K
00019530	D9 2A 88 EC 0F 14 15 A4 80 A2 02 2A E2 03 04 41	Û*^i...æ€¢.*â..A
00019540	84 3F A2 28 A2 0F C1 1D 71 65 11 9E B8 E0 82 A2	„?¢(¢.Á.qe.ž,à,¢
00019550	22 2A EB 63 69 81 96 AE B4 74 4B 5B BA AF D2 9D	"*ëci.-@'tK[°-Ò.
00019560	EE 74 4F D2 24 26 F9 F4 7F 66 6E 32 4D 26 77 32	îtOò\$&ùô.fn2M&w2
00019570	63 34 B5 40 CF 57 6D 33 73 EF B6 DB 20 F7 CF	u4wQžw +¢æxîô¢:ž

[그림 13] obj 문자를 기준으로 검색한 결과 - 1

<sup>2</sup> [https://github.com/py-pdf/pypdf/blob/82eac7e316f8f785d00ed600f8ba4aba3296a4a8/pypdf/\\_reader.py#L614](https://github.com/py-pdf/pypdf/blob/82eac7e316f8f785d00ed600f8ba4aba3296a4a8/pypdf/_reader.py#L614)

