

< [Go to the original](#)



Top 5 Open-Source PDF Scraping Tools Using Python

PDF scraping is often tricky since PDFs aren't designed for easy data extraction. Luckily, Python has several open-source tools that can...



Azhar Sayyad

Follow

a11y-light · October 25, 2024 (Updated: October 25, 2024) · **Free:** No

PDF scraping is often tricky since PDFs aren't designed for easy data extraction. Luckily, Python has several open-source tools that can help you extract data efficiently. Here's a look at the top 5 open-source PDF scraping libraries in Python, along with practical examples to get you started.

1. PyMuPDF (fitz)

- **Description:** PyMuPDF allows you to extract text, images, and metadata from PDFs. It works with both text-based and scanned

Freedium

- **Best for:** Extracting structured text and analyzing complex layouts.

- **Installation:**

Copy

```
pip install PyMuPDF
```

- **Example Usage:**

Copy

```
import fitz

doc = fitz.open("sample.pdf")
for page in doc:
    text = page.get_text()
    print(text)
doc.close()
```

- **Why Use It:** Fast, lightweight, and good for handling non-standard layouts.
- **Limitations:** Image extraction is basic.

2. pdfminer.six

- **Description:** An improved version of pdfminer, this library focuses on text extraction and layout analysis. It offers a lot of control when scraping PDFs with complex text.
- **Best for:** Extracting text with precise layout preservation.
- **Installation:**

Copy

- **Example Usage:**

Copy

```
from pdfminer.high_level import extract_text

text = extract_text('sample.pdf')
print(text)
```

- **Why Use It:** Ideal for PDFs with detailed layout and font information.
- **Limitations:** A bit slower compared to other tools.

3. pdfrw

- **Description:** pdfrw allows you to read and manipulate PDF files. It's great for merging, rotating, and extracting information from PDFs.
- **Best for:** Manipulating PDFs along with scraping.
- **Installation:**

Copy

```
pip install pdfrw
```

- **Example Usage:**

Copy

```
from pdfrw import PdfReader

pdf = PdfReader("sample.pdf")
```

- **Why Use It:** Simple and effective for both extraction and manipulation.
- **Limitations:** Limited support for text extraction compared to pdfminer.six.

4. PyPDF2

- **Description:** PyPDF2 is a popular library that helps with splitting, merging, and extracting text from PDFs. It works well with basic PDFs.
- **Best for:** Simple text extraction and PDF manipulation tasks.
- **Installation:**

Copy

```
pip install PyPDF2
```

- **Example Usage:**

Copy

```
from PyPDF2 import PdfReader

reader = PdfReader("sample.pdf")
for page in reader.pages:
    print(page.extract_text())
```

- **Why Use It:** Easy to use for basic tasks like splitting and merging PDFs.
- **Limitations:** Struggles with complex PDFs and doesn't support scanned PDFs.

- **Description:** camelot-py is a specialized library for extracting tables from PDFs. If you need to work with tabular data, this is the go-to tool.
- **Best for:** Extracting tables from PDFs accurately.
- **Installation:**

Copy

```
pip install camelot-py
```

- **Example Usage:**

Copy

```
import camelot

tables = camelot.read_pdf("sample.pdf", pages='1')
print(tables[0].df)
```

- **Why Use It:** Perfect for scraping tabular data from PDFs.
- **Limitations:** Works best with text-based PDFs, not scanned ones.

Summary Table

Library	Best For	Strengths	Limitations
PyMuPDF	Complex layouts	Fast, lightweight	Basic image extraction
pdfminer.six	Detailed layout preservation	Precise text extraction	Slower than others
pdfrw	Manipulating PDFs	Simple and effective	Limited text extraction
PyPDF2	Basic PDF tasks	Easy to use	Struggles with complex PDFs
camelot-py	Extracting tables	Great for tabular data	Not ideal for scanned PDFs

Choosing the right PDF scraping tool depends on your specific needs. If you're looking to extract tables, **camelot-py** is unbeatable. For more complex layouts, **PyMuPDF** or **pdfminer.six** would be your best bet. If your goal is simple PDF manipulation, **PyPDF2** or **pdfcrowd** should do the job.

These open-source tools make it much easier to scrape PDFs in Python, no matter how complex the document structure is.

[#python](#) [#pdf](#) [#pdf-scraping](#) [#extract](#) [#open-source](#)