

# Reading the PDF properties/metadata in Python

[Ask Question](#)

How can I read the properties/metadata like Title, Author, Subject and Keywords stored on a PDF file using Python?

[python](#)[pdf](#)[metadata](#)

edited Jun 2 at 20:10



wolphi

4,874 13 39

asked Jan 8 '13 at 6:13



Khaleel

1,301 3 16 33

## 4 Answers

Try [pdfminer](#):

```
from pdfminer.pdfparser
from pdfminer.pdfdocument

fp = open('diveintopython.pdf', 'rb')
parser = PDFParser(fp)
doc = PDFDocument(parser)

print doc.info # The "
```

Here's the output:

```
>>> [{'CreationDate': 'D:20130108061300+0000',
      'Creator': 'DocBook XSL Stylesheet Processor (v1.38.2)',
      'Keywords': 'Python, documentation, book, fr',
      'Producer': 'html2pdf',
      'Rights Reserved.',
      'Title': 'Dive Into Python'}]
```

For more info, look at this tutorial: [A lightweight XMP parser for extracting PDF metadata](#)

answered Jan 8 '13 at 6:22



**namit**

**4,505** 2 22 36

1 A heads-up: the author of pdfminer says it's incompatible with Python 3, at least as of date of this post ([link](#)) – [JSmyth](#) Jan 5 '14 at 22:36

8 As of November 2013, the "PDFDocument class now takes a PDFParser object as an argument. PDFDocument.set\_parser() and PDFParser.set\_document() is removed." So you can just do doc=PDFDocument(parser), and skip the calls to set\_document, set\_parser, and initialize. – [Derek Kurth](#) Oct 14 '14 at 15:55

@JSmyth The [PyPi Index](#) currently lists three working pdfminer forks that are compatible with Python 3. `pip search pdfminer` – [zero2cx](#) Jan 19 '17 at 7:10

@zero2cx thanks for the update. I personally settled on [pdfminer3k](#). Works well for my purposes. One has to read the API document in the repo though as the accepted here answer is not a valid API for pdfminer3k anymore. – [JSmyth](#) Jan 21 '17 at 23:33

1 There is now an official Python 3 fork of the project [github.com/pdfminer/pdfminer.six](https://github.com/pdfminer/pdfminer.six) – [Harsh](#) Jun 4 '17 at 12:46

For Python 3 see  
[PyPDF2](#) with example  
 code from @Khaleel  
 updated to:

```
from PyPDF2 import PdfFileReader
pdf_toread = PdfFileReader(pdf_file)
pdf_info = pdf_toread.getPage(0).get('/Info')
print(str(pdf_info))
```

Install using `pip install PyPDF2`.

answered Oct 8 '16 at 11:31



[Morten Zilmer](#)

12.3k 2 16 37

I have implemented this  
 using [pyPdf](#). Please see  
 the sample code below.

```
from pyPdf import PdfFileReader
pdf_toread = PdfFileReader(pdf_file)
pdf_info = pdf_toread.getPage(0).get('/Info')
print str(pdf_info)
```

Output:

```
{'/Title': u'Microsoft Word Document', '/CreationDate': u'D:2010.0.0 (Windows)', '/ModDate': u'D:2010.0.0 (Windows)', '/PScript5.dll Version': u'5.0'}
```

Note: pyPdf [homepage](#)  
 says it is no longer  
 maintained.

edited Jan 22 at 14:22

answered Jan 8 '13 at 8:49



[Khaleel](#)

1,301 3 16 33

1 Don't use `file`, use `open` instead. –  
[Burhan Khalid](#) Jan 9 '13 at 6:21

2 Note that the pyPdf is  
 marked on the

For Python 3 and new  
pdfminer (pip install  
pdfminer3k):

```
import os
from pdfminer.pdfparser
from pdfminer.pdfparser

fp = open("foo.pdf", 'r
parser = PDFParser(fp)
doc = PDFDocument(parse
parser.set_document(doc
doc.set_parser(parser)
if len(doc.info) > 0:
    info = doc.info[0]
    print(info)
```

answered Dec 19 '16 at 1:36



[Rabash](#)

695 8 10