Главная » Библиотеки Python

16.04.2021 \bigcirc 0

Библиотека PyPDF2 в Python

Программирование на Python

PyPDF2 – это библиотека Python для работы с файлами PDF. Мы можем использовать модуль PyPDF2 для работы с существующими файлами PDF, но мы не можем создать новый PDF-файл с помощью этого модуля.

```
Содержание ^
1. Возможности PyPDF2
2. Установка модуля
3. Примеры
4. 1. Извлечение метаданных PDF
5. 2. Извлечение текста из PDF-страниц
6. 3. Поворот страниц файла PDF
7. 4. Объединение файлов PDF
8. 5. Разделение PDF-файлов на одностраничные файлы
9. 6. Извлечение изображений
```

Некоторые из интересных особенностей модуля PyPDF2:

Возможности PyPDF2

Извлечение содержимого файла PDF постранично.

создатель, время создания и последнего обновления.

Метаданные PDF-файлов, такие как количество страниц, автор,

- Объедините несколько файлов PDF.
- Поворачивайте страницы PDF-файла на угол.

изображения с помощью библиотеки Pillow.

- Масштабирование страниц PDF. Извлечение изображений со страниц PDF и сохранение, как
- Установка модуля

Мы можем использовать PIP для установки модуля PyPDF2.

\$ pip install PyPDF2

Примеры

Рассмотрим несколько примеров работы с файлами PDF с помощью модуля PyPDF2.

1. Извлечение метаданных PDF

получить информацию об авторе PDF, приложении для создателя и датах создания.

import PyPDF2 3. with open('Python_Tutorial.pdf', 'rb') as pdf_file:

pdf reader = PyPDF2.PdfFileReader(pdf file)

Мы можем получить количество страниц в файле PDF. Мы также можем

```
print(f'Number of Pages in PDF File is
      {pdf_reader.getNumPages()}')
          print(f'PDF Metadata is {pdf_reader.documentInfo}')
          print(f'PDF File Author is {pdf_reader.documentInfo["/Author"]}')
          print(f'PDF File Creator is
      {pdf_reader.documentInfo["/Creator"]}')
Пример вывода:
      Number of Pages in PDF File is 2
      PDF Metadata is {'/Author': 'Microsoft Office User', '/Creator':
```

'Microsoft Word', '/CreationDate': "D:20191009091859+00'00'",

Файл PDF должен быть открыт в двоичном режиме. Вот почему

Класс PdfFileReader используется для чтения файла PDF. DocumentInfo – это словарь, содержащий метаданные файла PDF.

'/ModDate': "D:20191009091859+00'00'"}

PDF File Creator is Microsoft Word

атрибут numPages.

import PyPDF2

PDF File Author is Microsoft Office User

режим открытия файла передается, как 'rb'.

- Мы можем получить количество страниц в файле PDF с помощью функции getNumPages(). Альтернативный способ – использовать
- 2. Извлечение текста из PDF-страниц

with open('Python_Tutorial.pdf', 'rb') as pdf_file: pdf_reader = PyPDF2.PdfFileReader(pdf_file)

```
pdf_page = pdf_reader.getPage(0)
      print(pdf_page.extractText())
      for page_num in range(pdf_reader.numPages):
          pdf_page = pdf_reader.getPage(page_num)
          print(pdf_page.extractText())
Mетод getPage (int) PdfFileReader возвращает экземпляр
PyPDF2.pdf.PageObject.
Мы можем вызвать метод extractText() для объекта страницы, чтобы
```

ExtractText() не возвращает двоичных данных, таких как изображения.

получить текстовое содержимое страницы.

- 3. Поворот страниц файла PDF
- PyPDF2 допускает множество типов манипуляций, которые могут выполняться постранично. Мы можем повернуть страницу по часовой

import PyPDF2

стрелке или против часовой стрелки на угол.

with open('Python_Tutorial.pdf', 'rb') as pdf_file: pdf reader = PyPDF2.PdfFileReader(pdf_file) pdf_writer = PyPDF2.PdfFileWriter() for page_num in range(pdf_reader.numPages): pdf_page = pdf_reader.getPage(page_num)

pdf page.rotateClockwise(90) # rotateCounterClockwise()

```
pdf_writer.addPage(pdf_page)
     with open('Python_Tutorial_rotated.pdf', 'wb') as
 pdf_file_rotated:
         pdf_writer.write(pdf_file_rotated)
PdfFileWriter используется для записи файла PDF из исходного PDF.
Мы используем метод rotateClockwise (90), чтобы повернуть
страницу по часовой стрелке на 90 градусов.
Мы добавляем повернутые страницы в экземпляр PdfFileWriter.
Наконец, для создания повернутого файла PDF используется метод
```

write() класса PdfFileWriter. 4. Объединение файлов PDF

import PyPDF2

выпуске GitHub.

import contextlib

import PyPDF2

pdf_files_list]

for f in files:

записи.

pdf_merger = PyPDF2.PdfFileMerger()

pdf_files_list = ['Python_Tutorial.pdf',

pdf_merger.write(pdf_file_merged)

'Python_Tutorial_rotated.pdf'] for pdf_file_name in pdf_files_list: with open(pdf_file_name, 'rb') as pdf_file: pdf_merger.append(pdf_file)

with open('Python_Tutorial_merged.pdf', 'wb') as pdf_file_merged:

Приведенный выше код подходит для объединения файлов PDF, но он создал пустой PDF-файл. Причина в том, что исходные файлы PDF были закрыты до того, как произошла фактическая запись для создания объединенного файла PDF.

Это ошибка последней версии PyPDF2. Вы можете прочитать об этом в

Существует альтернативный подход к использованию модуля contextlib,

чтобы исходные файлы оставались открытыми до завершения операции

4. pdf_files_list = ['Python_Tutorial.pdf', 'Python_Tutorial_rotated.pdf'] 6. with contextlib.ExitStack() as stack:

files = [stack.enter_context(open(pdf, 'rb')) for pdf in

with open('Python_Tutorial_merged_contextlib.pdf', 'wb') as f:

pdf_merger = PyPDF2.PdfFileMerger()

pdf_merger.append(f)

pdf_merger.write(f)

5. Разделение PDF-файлов на

Python_Tutorial_0.pdf и Python_Tutorial_1.pdf.

6. Извлечение изображений

следующую команду.

\$ pip install Pillow

изображения из файла PDF.

```
одностраничные файлы
      import PyPDF2
     with open('Python_Tutorial.pdf', 'rb') as pdf_file:
         pdf_reader = PyPDF2.PdfFileReader(pdf_file)
         for i in range(pdf_reader.numPages):
             pdf_writer = PyPDF2.PdfFileWriter()
             pdf_writer.addPage(pdf_reader.getPage(i))
             output_file_name = f'Python_Tutorial_{i}.pdf'
             with open(output_file_name, 'wb') as output_file:
                 pdf_writer.write(output_file)
Python_Tutorial.pdf состоит из 2 страниц. Выходные файлы называются
```

Мы можем использовать PyPDF2 вместе с Pillow (Python Imaging Library)

Прежде всего, вам нужно будет установить модуль Pillow, используя

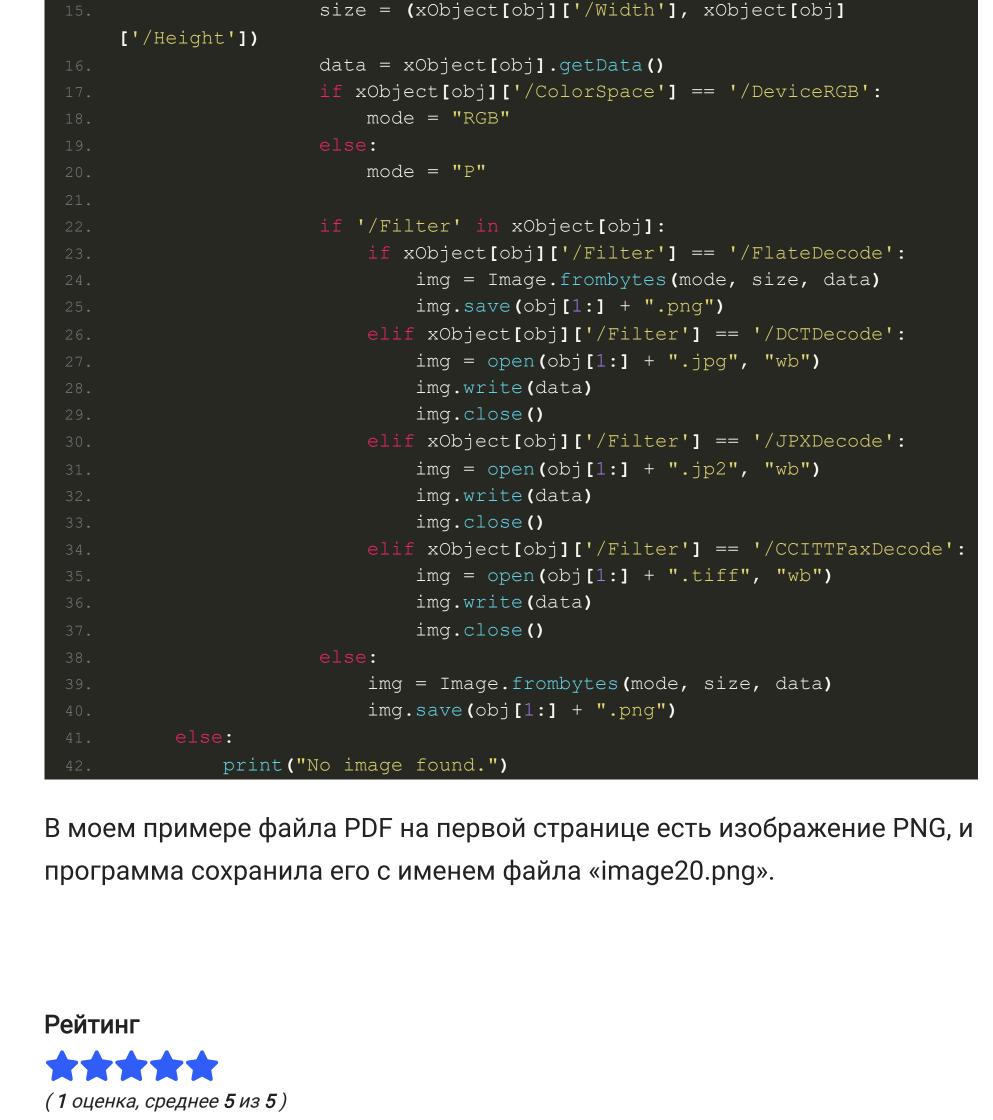
Вот простая программа для извлечения изображений с первой страницы

файла PDF. Мы можем легко расширить его, чтобы извлечь все

для извлечения изображений со страниц PDF и сохранения их в виде файлов изображений.

4. with open('Python_Tutorial.pdf', 'rb') as pdf_file: pdf_reader = PyPDF2.PdfFileReader(pdf_file)

page0 = pdf_reader.getPage(0) if '/XObject' in page0['/Resources']: xObject = page0['/Resources']['/XObject'].getObject() for obj in xObject: if xObject[obj]['/Subtype'] == '/Image':



Помогаю в изучении Питона на примерах. Автор практических задач с

Похожие материалы

Библиотека SciKit Learn в

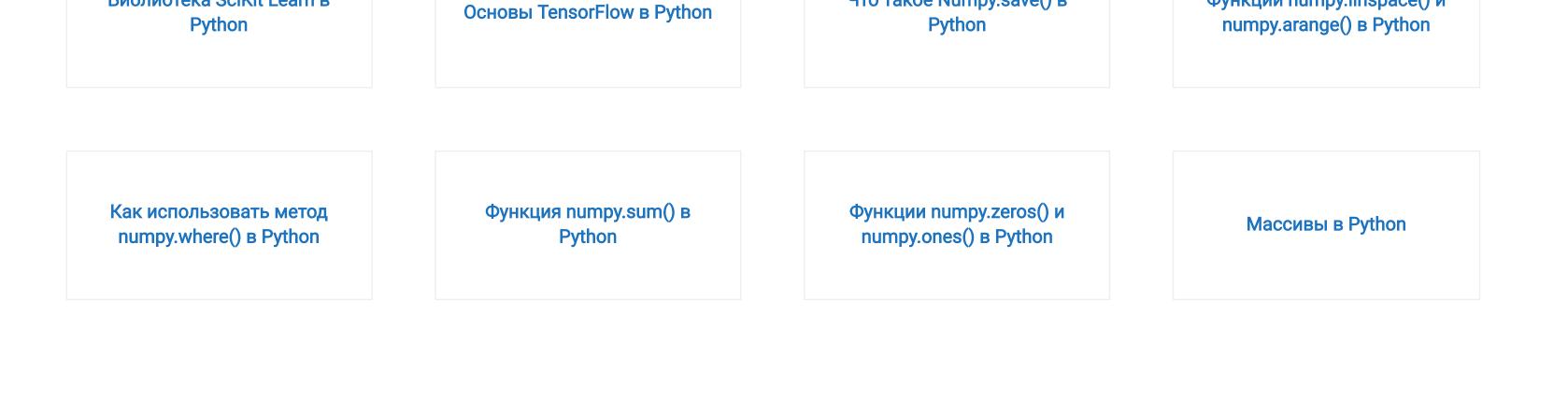
© 2022 Блог Васильева А.Н.

Комментарии ⁰

Васильев А.Н. / автор статьи

детальным разбором их решений.

Поделиться:



Что такое Numpy.save() в

Функции numpy.linspace() и