

---

# **Red Stork**

***Release 0.0.41***

**Jul 09, 2020**



|          |                            |           |
|----------|----------------------------|-----------|
| <b>1</b> | <b>Quick Start</b>         | <b>3</b>  |
| 1.1      | Tutorial . . . . .         | 3         |
| 1.2      | API Reference . . . . .    | 6         |
| <b>2</b> | <b>Indices and tables</b>  | <b>13</b> |
|          | <b>Python Module Index</b> | <b>15</b> |
|          | <b>Index</b>               | <b>17</b> |



Yet another PDF parser. This one is based on [PDFium](#) engine.



# CHAPTER 1

---

## Quick Start

---

Sample:

```
from redstork.document import Document

doc = Document('sample.pdf')
print('Number of pages:', len(doc))
```

## 1.1 Tutorial

The philosophy of **redstork** is to map API to standard and well understood Python objects, like `list` and `dict`. In this tutorial we will use the following sample document.

### 1.1.1 Version

There are two version values in `redstork` module: PDFium build version, and Python package version:

```
import redstork

redstork.__pdfium_version__
>> 'cromium/4097'

redstork.__version__
>> '0.0.1'
```

### 1.1.2 Document

`Document` is the top-level object, and the only object that can be instantiated directly:

```
from redstork import Document

doc = Document('sample.pdf')

len(doc)
>> 15
```

As you can see, *Document* resembles standard Python *list*, containing *Page* objects.

PDF file creators can attach arbitrary key-value strings to the document, that we call meta (official PDF specs call it Document Information Dictionary). Most commonly these values describe Author, Title, and the name of software that created this document. Lets see the meta in our sample:

```
doc.meta['Title']
>> 'Red Stork'
```

You can change meta content and save the updated document:

```
doc.meta['Title'] = 'Awesome PDF parsing library'
doc.save('awesome.pdf')
```

Document has a lazily populated collection of fonts. Initially this collection is empty. As pages are being accessed and parsed, this collection is being populated:

```
list(doc[0]) # read all objects from page 1
len(doc.fonts)
>> 2
```

### 1.1.3 Page

*Page* represents PDF page. Get page by indexing a *Document* object, just like a normal list:

```
page = doc[0]
page.crop_box
>> (0.0, 0.0, 612.0, 792.0)
```

*Page* has *Page.label*, representing the page label (like xxi, or 128):

```
doc[2].label # this is the label of the third page
>> 'i'
```

A page of PDF document is a list-like object, containing concrete instances of *PageObject*:

```
page = ...
len(page) # how many objects on this page?
>> 17
```

### 1.1.4 PageObject

Abstract class *PageObject* describes an object on a PDF page. Concrete classes implementing *PageObject* are:

- *TextObject* - a string of characters
- *PathObject* - vector graphics
- *ImageObject* - a bitmap image



- *ShadingObject* - a shading object

Notable properties of all objects are:

- *PageObject.page()* - links back to the parent page
- *PageObject.matrix()* - transformation matrix of this object
- *PageObject.rect()* - rectangle of this object on the page

### 1.1.5 TextObject

Text object represents a string of characters. Each character is a three-tuple of (charcode, x, y), where charcode is a character code (this value is just an index in the font glyph table, not a text corresponding to this character!). x and y are placement coordinates of this character (in the coordinate system of this *TextObject* - first character typically has x, y == 0, 0).

Text object has font property. Here is how to use font to extract text of a *TextObject*:

```
def text_of(o):
    assert o.type == PageObject.OBJ_TYPE_TEXT, o
    text = []
    for c, x, y in o:
        text.append(o.font[c])
    return ''.join(text)

page = ...
for o in page:
    if o.type == PageObject.OBJ_TYPE_TEXT:
        text = text_of(o)
        print(text)
```

### 1.1.6 PathObject

Path object represents a set of vector drawing instructions.

### 1.1.7 ImageObject

Image object represents an embedded bitmap image. You can get the pixel width and height of the image, using the properties *ImageObject.pixel\_width()* and *ImageObject.pixel\_height()*.

Example:

```
page = ...
for o in page:
    if o.type == PageObject.OBJ_TYPE_IMAGE:
        print(o.pixel_width, o.pixel_height)
```

### 1.1.8 Font

Font object is a look-up table for character text, and also holds character glyphs (shape).

Font names in PDF file have a special prefix. To get a human-friendly one use *Font.simple\_name()*.

Document contains a lazy font collection `Document.fonts()`. It is lazy, because just after document is opened, it is empty. As pages are accessed and parsed, this collection is populated.

Here is how to get `Glyph` object:

```
page = ...
for o in page:
    if o.type == PageObject.OBJ_TYPE_TEXT:
        for code, _, _ in o:
            glyph = o.font.load_glyph(code)
            print('Character with code %d has %d glyph instructions', code, len(glyph))
```

## 1.2 API Reference

### 1.2.1 Document

**class** `redstork.Document` (*file\_name*, *password=None*)  
PDF document.

A list-like container of pages. Sample use:

```
doc = Document('sample.pdf')
print("Number of pages:", len(doc))

for key, value in doc.meta.items():
    print(' ', key, ': ', value)
```

**\_\_init\_\_** (*file\_name*, *password=None*)  
Create new PDF Document object, from a file.

#### Parameters

- **file\_name** (*str*) – Name of PDF file
- **password** (*str*) – File password (optional)

**numpages** = `None`  
*int* – total number of pages

**meta** = `None`  
*dict* – document meta info (Author, Title, etc)

**fonts** = `None`  
*dict* – font collection (populated lazily as pages are parsed)

**\_\_getitem\_\_** (*page\_index*)  
Returns `Page` at this index.

Example:

```
doc = ...

page = doc[0] # first page
```

**Parameters** **page\_index** (*int*) – zero-based page index

**Returns** `Page` object

**\_\_len\_\_()**  
Returns number of pages in this document

**\_\_iter\_\_()**  
Iterate over the pages of this document

**get\_all\_pages\_size()**  
get width and height of all pages, without loading each page

**changed**  
True is PDF was changed since the load (or last save)

**save(filename)**  
Saves PDF file, resets `Document.changed()` to False

### 1.2.2 Page

**class redstork.Page**(page, page\_index, parent)  
Represents page of a PDF file.

**crop\_box**  
Page crop box.

**media\_box**  
Page media box.

**rotation**  
Page rotation.

- 0 - no rotation
- 1 - rotated 90 degrees clock-wise
- 2 - rotated 180 degrees clock-wise
- 3 - rotated 270 degrees clock-wise

**label**  
Page label.

**\_\_len\_\_()**  
Number of objects on this page.

**\_\_getitem\_\_**(index)  
Get object at this index.

**\_\_iter\_\_()**  
Iterates over page objects.

**flat\_iter()**  
Iterates over all non-container objects (Text, Image, Path).

**render\_to\_buffer**(scale=1.0, rect=None)  
Render page (or rectangle on the page) to memory (the pixel format is BGRx)

#### Parameters

- **scale** (*float*) – scale to use (default is 1.0, which will assume that 1pt takes 1px)
- **rect** (*tuple*) – optional rectangle to render. Value is a 4-tuple of (x0, y0, x1, y1) in PDF coordinates. if None, then page's `crop_box` will be used for rendering.

**render**(file\_name, scale=1.0, rect=None)  
Render page (or rectangle on the page) as PPM image file.

### Parameters

- **file\_name** (*str*) – name of the output file
- **scale** (*float*) – scale to use (default is 1.0, which will assume that 1pt takes 1px)
- **rect** (*tuple*) – optional rectangle to render. Value is a 4-tuple of (x0, y0, x1, y1) in PDF coordinates. if None, then page's *crop\_box* will be used for rendering.

## PageObject

**class** redstork.**PageObject** (*obj, index, typ, parent*)

**OBJ\_TYPE\_TEXT** = 1  
see *TextObject*

**OBJ\_TYPE\_PATH** = 2  
see *PathObject*

**OBJ\_TYPE\_IMAGE** = 3  
see *ImageObject*

**OBJ\_TYPE\_SHADING** = 4  
see *ShadingObject*

**OBJ\_TYPE\_FORM** = 5  
Common superclass of all page objects

**type** = None  
type of this object

**matrix** = None  
transformation matrix of this object

**page**  
Links back to the parent page

## TextObject

**class** redstork.**TextObject** (*obj, index, typ, parent*)

Represents a string of text on a page

**font** = None  
*Font* for this text object

**font\_size** = None  
font size of this text object

**matrix** = None  
matrix for this page object

**\_\_len\_\_**()  
Number of items in this string

**\_\_getitem\_\_** (*index*)  
Returns item at this index.  
Each item is a 3-tuple: (charcode, x, y).

**\_\_iter\_\_**()  
Iterates over items.

**char\_iter()**  
Iterates over characters (skips kerns)

**text\_geometry\_iter()**  
Iterates over characters and returns character text and bounds

**effective\_font\_size**  
Returns effective (user-visible) font size

**scale\_y**  
Returns Y-scale of text matrix transformation

**scale\_x**  
Returns X-scale of text matrix transformation

**skew**  
Returns skew value of text matrix.

**box** (*x0*, *y0*, *x1*, *y1*)  
Computes bounding box after transformation with text matrix

## Font

**class** `redstork.Font` (*font*, *parent*)  
Represents font used in a PDF file.

**FLAGS\_NORMAL** = 0  
Normal font

**FLAGS\_FIXED\_PITCH** = 1  
Fixed pitch font

**FLAGS\_SERIF** = 2  
Serif font

**FLAGS\_SYMBOLIC** = 4  
Symbolic font

**FLAGS\_SCRIPT** = 8  
Script font

**FLAGS\_NONSYMBOLIC** = 32  
Non-symbolic font

**FLAGS\_ITALIC** = 64  
Italic font

**FLAGS\_ALLCAP** = 65536  
All-cap font

**FLAGS\_SMALLCAP** = 131072  
Small-cap font

**FLAGS\_FORCE\_BOLD** = 262144  
Force-bold font

**name**  
Font name in the PDF document.

**simple\_name**  
Font name without PDF-specific prefix.

**flags**

Font flags.

**weight**

Font weight.

**is\_vertical**

True for vertical writing systems (CJK)

**id**

Tuple of (Object\_id, Generation\_id), identifying underlying stream in PDF file

**load\_glyph** (*charcode*)

Load glyph, see *Glyph*

**Parameters** **charcode** (*int*) – the character code (see *TextObject*)

**\_\_getitem\_\_** (*charcode*)

Returns Unicode text of this character.

**Parameters** **charcode** (*int*) – the character code (see *TextObject*) –

**is\_editable**

True if font encoding can be changed

**\_\_setitem\_\_** (*charcode, text*)

Updates font encoding.

**Parameters**

- **charcode** (*int*) – character code
- **text** (*str*) – new text for this character code

**Raises** *ReadOnlyEncodingError* – if encoding is read-one (no “ToUnicode” map in the font dictionary)

## Glyph

**class** *redstork.Glyph* (*glyph, parent*)

Represents Glyph drawing instructions

**LINETO = 0**

LineTo instruction

**CURVETO = 1**

CurveTo instruction

**MOVETO = 2**

MoveTo instruction

**\_\_getitem\_\_** (*i*)

Returns a 4-tuple representing this drawing instruction: (x, y, type, close).

**Parameters** **i** (*int*) – index of the instruction

## ImageObject

**class** *redstork.ImageObject* (*obj, index, typ, parent*)

Represents image on a page.

**matrix = None**  
matrix for this page object

**pixel\_width**  
width of the bitmap, in pixels

**pixel\_height**  
height of the bitmap, in pixels

### PathObject

**class** redstork.**PathObject** (*obj, index, typ, parent*)  
Represents vector graphics on a aage.

**matrix = None**  
matrix for this page object

### ShadingObject

**class** redstork.**ShadingObject** (*obj, index, typ, parent*)  
Represents a shading object on a page.

### FormObject

**class** redstork.**FormObject** (*obj, index, typ, parent*)  
Represents a form (XObject) on a page - a container of other page objects (used internally).

**matrix = None**  
matrix for this page object

**form\_matrix = None**  
transformation matrix for contained objects

**flat\_iter()**  
Iterates over all non-container objects in this form.





## CHAPTER 2

---

### Indices and tables

---

- [genindex](#)
- [modindex](#)
- [search](#)
- [Glossary](#)



**r**

redstork, [6](#)



## Symbols

[\\_\\_getitem\\_\\_\(\)](#) (*redstork.Document* method), 6  
[\\_\\_getitem\\_\\_\(\)](#) (*redstork.Font* method), 10  
[\\_\\_getitem\\_\\_\(\)](#) (*redstork.Glyph* method), 10  
[\\_\\_getitem\\_\\_\(\)](#) (*redstork.Page* method), 7  
[\\_\\_getitem\\_\\_\(\)](#) (*redstork.TextObject* method), 8  
[\\_\\_init\\_\\_\(\)](#) (*redstork.Document* method), 6  
[\\_\\_iter\\_\\_\(\)](#) (*redstork.Document* method), 7  
[\\_\\_iter\\_\\_\(\)](#) (*redstork.Page* method), 7  
[\\_\\_iter\\_\\_\(\)](#) (*redstork.TextObject* method), 8  
[\\_\\_len\\_\\_\(\)](#) (*redstork.Document* method), 6  
[\\_\\_len\\_\\_\(\)](#) (*redstork.Page* method), 7  
[\\_\\_len\\_\\_\(\)](#) (*redstork.TextObject* method), 8  
[\\_\\_setitem\\_\\_\(\)](#) (*redstork.Font* method), 10

## B

[box\(\)](#) (*redstork.TextObject* method), 9

## C

[changed](#) (*redstork.Document* attribute), 7  
[char\\_iter\(\)](#) (*redstork.TextObject* method), 9  
[crop\\_box](#) (*redstork.Page* attribute), 7  
[CURVETO](#) (*redstork.Glyph* attribute), 10

## D

[Document](#) (*class in redstork*), 6

## E

[effective\\_font\\_size](#) (*redstork.TextObject* attribute), 9

## F

[flags](#) (*redstork.Font* attribute), 9  
[FLAGS\\_ALLCAP](#) (*redstork.Font* attribute), 9  
[FLAGS\\_FIXED\\_PITCH](#) (*redstork.Font* attribute), 9  
[FLAGS\\_FORCE\\_BOLD](#) (*redstork.Font* attribute), 9  
[FLAGS\\_ITALIC](#) (*redstork.Font* attribute), 9  
[FLAGS\\_NONSYMBOLIC](#) (*redstork.Font* attribute), 9  
[FLAGS\\_NORMAL](#) (*redstork.Font* attribute), 9

[FLAGS\\_SCRIPT](#) (*redstork.Font* attribute), 9  
[FLAGS\\_SERIF](#) (*redstork.Font* attribute), 9  
[FLAGS\\_SMALLCAP](#) (*redstork.Font* attribute), 9  
[FLAGS\\_SYMBOLIC](#) (*redstork.Font* attribute), 9  
[flat\\_iter\(\)](#) (*redstork.FormObject* method), 11  
[flat\\_iter\(\)](#) (*redstork.Page* method), 7  
[Font](#) (*class in redstork*), 9  
[font](#) (*redstork.TextObject* attribute), 8  
[font\\_size](#) (*redstork.TextObject* attribute), 8  
[fonts](#) (*redstork.Document* attribute), 6  
[form\\_matrix](#) (*redstork.FormObject* attribute), 11  
[FormObject](#) (*class in redstork*), 11

## G

[get\\_all\\_pages\\_size\(\)](#) (*redstork.Document* method), 7  
[Glyph](#) (*class in redstork*), 10

## I

[id](#) (*redstork.Font* attribute), 10  
[ImageObject](#) (*class in redstork*), 10  
[is\\_editable](#) (*redstork.Font* attribute), 10  
[is\\_vertical](#) (*redstork.Font* attribute), 10

## L

[label](#) (*redstork.Page* attribute), 7  
[LINETO](#) (*redstork.Glyph* attribute), 10  
[load\\_glyph\(\)](#) (*redstork.Font* method), 10

## M

[matrix](#) (*redstork.FormObject* attribute), 11  
[matrix](#) (*redstork.ImageObject* attribute), 10  
[matrix](#) (*redstork.PageObject* attribute), 8  
[matrix](#) (*redstork.PathObject* attribute), 11  
[matrix](#) (*redstork.TextObject* attribute), 8  
[media\\_box](#) (*redstork.Page* attribute), 7  
[meta](#) (*redstork.Document* attribute), 6  
[MOVETO](#) (*redstork.Glyph* attribute), 10

## N

`name` (*redstork.Font* attribute), 9  
`numpages` (*redstork.Document* attribute), 6

## O

`OBJ_TYPE_FORM` (*redstork.PageObject* attribute), 8  
`OBJ_TYPE_IMAGE` (*redstork.PageObject* attribute), 8  
`OBJ_TYPE_PATH` (*redstork.PageObject* attribute), 8  
`OBJ_TYPE_SHADING` (*redstork.PageObject* attribute), 8  
`OBJ_TYPE_TEXT` (*redstork.PageObject* attribute), 8

## P

`Page` (class in *redstork*), 7  
`page` (*redstork.PageObject* attribute), 8  
`PageObject` (class in *redstork*), 8  
`PathObject` (class in *redstork*), 11  
`pixel_height` (*redstork.ImageObject* attribute), 11  
`pixel_width` (*redstork.ImageObject* attribute), 11

## R

`redstork` (module), 6  
`render()` (*redstork.Page* method), 7  
`render_to_buffer()` (*redstork.Page* method), 7  
`rotation` (*redstork.Page* attribute), 7

## S

`save()` (*redstork.Document* method), 7  
`scale_x` (*redstork.TextObject* attribute), 9  
`scale_y` (*redstork.TextObject* attribute), 9  
`ShadingObject` (class in *redstork*), 11  
`simple_name` (*redstork.Font* attribute), 9  
`skew` (*redstork.TextObject* attribute), 9

## T

`text_geometry_iter()` (*redstork.TextObject* method), 9  
`TextObject` (class in *redstork*), 8  
`type` (*redstork.PageObject* attribute), 8

## W

`weight` (*redstork.Font* attribute), 10