

PDF Text Extractors

Tools for Extracting Text from PDFs

Danny S. Vargas
Viviane Moreira

08/05/2019

- 1 Tika
- 2 PDFBox
- 3 mupdf
- 4 pdfminer
- 5 pdftotext
- 6 poppler(popplerqt4)
- 7 textract
- 8 pypdf2xml

- Baseline
- Java

- Tika
- Java

- PDF, XPS, and E-book viewer
- C

- Exact location of text in a page
- Python

- xpdf
- poppler

- PDF viewer
- xpdf3.0
- python

- Any type of file
- pdftotext and pdfminer
- python

- Python implementation of pdftoxml
- pdfminer(unicode characters)

Common Features

- PDFBox & Tika (Java)
- pdfminer, pdftotex & textract
- pdftotext & poppler(poppler-utils)
- pypdf2xml(pdftoxml in python)

- Baseline(Tika)

- No spaces between pages
- Similar results (structure and content)

- Additional break-lines
- Similar results(structure and content)

- Results depend on the input
 - (+) Follow the correct order
 - (-) One character per line
- Delayed

- Similar performance to Tika
- It does not follow the same order

- It maintains the structure (position of the text in the pdf document)

- pdftotext and pdfminer(methods)
- Results depend on the method used

Tika vs Others & pdf2xml

- Organized
- How to interpret and clean the tags?
- Used as reference of time performance

Summary

- Dependency on the input
- Correct order
- Configuration
- Documentation
- Time performance

Finally

- PDFBox
- Tika
- mupdf

- Github repository:
`https://github.com/dannysv/extract_textpdf`.