_/_/_

**PDF :-** Portable Document format. Generally PDF file Contains text, numbers, images, fonts, tables, special characters.

**Pdfminer :-**

→ It is used to extracting the text data from Pdf file.

→ It allow access to text position & font information.

→ It support various output formats like plain Text, HTML, XML.

**Pdf query :-**

→ It is built on top of the pdf miner.

→ It is designed for extracting structured data from Pdf file.

→ it is used to extract the tabular data.

**PyPdf2 :-** → Adding or removing pages.

→ It supports merging, splitting, rotation of PDF pages.

→ It allows text extraction & text search within PDF.

**Pymupdf :-**

→ It is the high performance Pdf library.

→ It supports text extraction & image extraction etc.

**PDF Parser :-**

→ It is a s/w Component.

→ It is used to extract information & content from PDF.

**Need for parsing :-**

To avoid errors and inaccurate data extraction we use parsing.

**Text extraction :-**

Primary task of PDF parser is to extract text from PDF document

## Font Handling:-

PDF file uses different fonts for text. PDF parser is responsible for handling font information.

## Image extraction:-

PDF parser is responsible for images, graphs extraction. extract images in various format & Resolution

## Metadata extraction:-

ex:- title, author, creation Date.

## Content Search:-

Search specific word with in the file.

## Document Manipulation:-

Adding, or removing pages, merging, splitting etc.

## Conversion:-

Converts PDF to plain text, HTML, XML.

## Program:-

```
from    nltk.tokenize import RegexpTokenizer
from    pdfminer.high_level import extract_text
from    nltk.probability import FreqDist.

# extract the text from PDF file
text = extract_text('/mahi/Dw/2010.00462.pdf')

# create an instance of tokenizer using NLTK
                    RegexpTokenizer
tokenizer = RegexpTokenizer('\w+')

# Tokenize the text read from PDF
    tokens = tokenizer.tokenize(text)
```

\w+ pattern only word characters are included in ~~tokens~~ Tokens (like letters, digits & underscores). removes non-word characters (.,!)

```
# find frequency Distribution
       freqdist = Freq.Dist (tokens)
# find words whose length is greater than 5
   and frequency greater than 20
long_frequent_words = [ words for words in tokens
                    if len(words) > 5 and
                    freqdist [words] > 20 ]

long_frequent_words
FreqDist (long_frequent_words) . plot ()
```

Google calab :-
```
       from google.calab import files
       files . upload ()
```