unit Test - 05

① Web scraping is a term used to describe the use of a program or algorithm to extract & process large amounts of data from the web.

Downloading a web page

```
>>> import requests
>>> res = requests. get ('http:// www. gutenberg. org/cach/
                         epub/1112/pg 111 2.txt.')
>>> type (res)
   <class 'requests. models. Response'>
>>> res. status - code == requests. codes. ok
   True
>>>len (res. text )
   178981
>>> print (res. text [: 250])
```

Saving Downloaded files to the Hard Drive

```
>>> import requests
>>> res = requests. get ('http: //www. gutenberg. org/cach/
                         Epub/1112 /pg "12. text ')
>>> res. raise - for - status ()
>>> playfile = open ('Romeo And Juliet. txt' ,'wb')
>>> for chunk in res. iter - content (100000):
        playfile. write (chunk)
   100000
   78981
>>> play file. close ()
```

(2) <u>Creating PDF's</u>

① Open one or more Existing PDF's into pdf file Reader objects.

② Create a new pdf file writer object.

③ Copy pages from the pdf file Reader objects into the Pdf File Writer object.

④ Finally, use the Pdf file Writer object to write the output PDF.

<u>Copying Pages</u>

```
import PyPDF2
pdf1 file = open('meeting minutes.pdf', 'rb')
pdf2 file = open('meeting minutes.pdf', 'rb')
pdf1 Reader = PyPDF2. Pdf File Reader (pdf1 File)
pdf2 Reader = PyPDF2. Pdf File Reader ( pdf2 File)
pdf Writer = PyPDF2. pdf File Writer ()

for pageNum insange ( pdf1 Reader, num Pages ):
    page Obj = pdf1 Reader. get page (pageNum)
    pdf Writer. add Page (pageObj)


for pageNum in range ( pdf2 Reader. num Pages ):
    pageObj = pdf2 Reader. get Page ( pageNum)
    pdf Writer. add Page (pageObj)


pdf Output file = open ('combined minutes. pdf', 'wb')
pdf Writer. write (pdf Output files)
Pdf Output File. close ()
Pdf1 File. close ()
pdf2 File. close ()
```

## Rotating Pages:

```
import PyPDF2
minutesfile = open('meeting minutes.pdf', 'rb')
pdfReader = PyPDF2. PdfFileReader (minutesFile)
page = pdfReader.getPage (0)
page. rotateClockwise (90)

pdfWriter = PyPDF2. PdfFileWriter ()
pdfWriter.addPage (page)
result pdffile = open ('rotated.Page.pdf', 'wb')
pdfWriter. write (result Pdf file)
result Pdf File. close ()
minutesFile. close ()
```