

```
# install tabula python package
!pip install tabula.py
```



```
Collecting tabula.py
  Downloading tabula_py-2.9.0-py3-none-any.whl (12.0 MB)
    12.0/12.0 MB 19.1 MB/s
Requirement already satisfied: pandas>=0.25.3 in /usr/local/lib/python
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist
Requirement already satisfied: distro in /usr/lib/python3/dist-package
Requirement already satisfied: python-dateutil>=2.8.1 in /usr/local/li
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/d
Installing collected packages: tabula.py
Successfully installed tabula.py-2.9.0
```

```
!pip install tabulate
```



```
➞ Requirement already satisfied: tabulate in /usr/local/lib/python3.10/d
```

```
# import the necessary libraries
from tabula import read_pdf
from tabulate import tabulate
```



```
import warnings
```

```
# ignore all warnings
warnings.filterwarnings("ignore")
```

```
# filename variable of the pdf file which needs to be uploaded into the fc
pdf_file = 'FoodList.pdf'
```

```
# extract data from page 1 of the pdf file
page_number = 1
```

```
# returns the extracted tables as pandas dataframes
tables_df = read_pdf(pdf_file, pages=page_number)
```

```
# print the tables from page 1 of the pdf
print(tables_df)
```

```
# ignore any warnings
```

```
WARNING:tabula.backend:Error importing jpye dependencies. Fallback to
WARNING:tabula.backend:No module named 'jpye'
WARNING:tabula.backend:Got stderr: Apr 01, 2024 4:58:47 PM org.apache.
WARNING: New fonts found, font cache will be re-built
Apr 01, 2024 4:58:47 PM org.apache.pdfbox.pdmodel.font.FileSystemFontP
WARNING: Building on-disk font cache, this may take a while
Apr 01, 2024 4:58:48 PM org.apache.pdfbox.pdmodel.font.FileSystemFontP
WARNING: Finished building on-disk font cache, found 17 fonts
```

	BREADS & CEREALS	Portion size * \
0	Bagel (1 average)	140 cal (45g)
1	Biscuit digestives	86 cal (per biscuit)
2	Jaffa cake	48 cal (per biscuit)
3	Bread white (thick slice)	96 cal (1 slice 40g)
4	Bread wholemeal (thick)	88 cal (1 slice 40g)
5	Chapatis	250 cal
6	Cornflakes	130 cal (35g)
7	Crackerbread	17 cal per slice
8	Cream crackers	35 cal (per cracker)
9	Crumpets	93 cal (per crumpet)
10	Flapjacks basic fruit mix	320 cal
11	Macaroni (boiled)	238 cal (250g)
12	Muesli	195 cal (50g)
13	Naan bread (normal)	300 cal (small plate size)
14	Noodles (boiled)	175 cal (250g)
15	Pasta (normal boiled)	330 cal (300g)
16	Pasta (wholemeal boiled)	315 cal (300g)
17	Porridge oats (with water)	193 cal (350g)
18	Potatoes** (boiled)	210 cal (300g)
19	Potatoes** (roast)	420 cal (300g)

	per 100 grams (3.5 oz)	Unnamed: 0	energy content
0	310 cal	NaN	Medium
1	480 cal	NaN	High
2	370 cal	NaN	Med-High
3	240 cal	NaN	Medium
4	220 cal	NaN	Low-med
5	300 cal	NaN	Medium
6	370 cal	NaN	Med-High
7	325 cal	NaN	Low Calorie
8	440 cal	NaN	Low / portion
9	198 cal	NaN	Low-Med
10	500 cal	NaN	High
11	95 cal	NaN	Low calorie
12	390 cal	NaN	Med-high
13	320 cal	NaN	Medium

14	70 cals	NaN	Low calorie
15	110 cals	NaN	Low calorie
16	105 cals	NaN	Low calorie
17	55 cals	NaN	Low calorie
18	70 cals	NaN	Low calorie
19	140 cals	NaN	Medium]

```
# use list comprehension to create a new list, loop through each dataframe
cleaned_tables = [table.dropna(axis='columns') for table in tables_df]

# Loop through the table and print everything, should not have any NaN val
for idx, table in enumerate(cleaned_tables) :
    print(f"Table {idx+1} after dropping NaN values:")

    print (table)
```

Table 1 after dropping NaN values:

BREADS & CEREALS		Portion size *	\
0	Bagel (1 average)	140 cal (45g)	
1	Biscuit digestives	86 cal (per biscuit)	
2	Jaffa cake	48 cal (per biscuit)	
3	Bread white (thick slice)	96 cal (1 slice 40g)	
4	Bread wholemeal (thick)	88 cal (1 slice 40g)	
5	Chapatis	250 cal	
6	Cornflakes	130 cal (35g)	
7	Crackerbread	17 cal per slice	
8	Cream crackers	35 cal (per cracker)	
9	Crumpets	93 cal (per crumpet)	
10	Flapjacks basic fruit mix	320 cal	
11	Macaroni (boiled)	238 cal (250g)	
12	Muesli	195 cal (50g)	
13	Naan bread (normal)	300 cal (small plate size)	
14	Noodles (boiled)	175 cal (250g)	
15	Pasta (normal boiled)	330 cal (300g)	
16	Pasta (wholemeal boiled)	315 cal (300g)	
17	Porridge oats (with water)	193 cal (350g)	
18	Potatoes** (boiled)	210 cal (300g)	
19	Potatoes** (roast)	420 cal (300g)	

per 100 grams (3.5 oz) energy content		
0	310 cal	Medium
1	480 cal	High
2	370 cal	Med-High
3	240 cal	Medium
4	220 cal	Low-med
5	300 cal	Medium
6	370 cal	Med-High
7	325 cal	Low Calorie
8	440 cal	Low / portion
9	198 cal	Low-Med
10	500 cal	High
11	95 cal	Low calorie
12	390 cal	Med-high
13	320 cal	Medium
14	70 cal	Low calorie
15	110 cal	Low calorie
16	105 cal	Low calorie
17	55 cal	Low calorie
18	70 cal	Low calorie
19	140 cal	Medium

```
# extract data from page 1 of the paf file
page_number = 3

# returns the extracted tables as pandas dataframes
tables_df = read_pdf(pdf_file, pages=page_number)

# print the tables from page 1 of the pdf
print(tables_df)
```

[Fish cake	90 cal	per cake	200 cal	Mediu
0	Fish fingers	50 cal	per piece	220 cal	Medium
1	Gammon	320 cal		280 cal	Med-High
2	Haddock fresh	200 cal		110 cal	Low calorie
3	Halibut fresh	220 cal		125 cal	Low calorie
4	NaN	NaN		NaN	NaN
5	Ham	6 cal		240 cal	Medium
6	Herring fresh grilled	300 cal		200 cal	Medium
7	Kidney	200 cal		160 cal	Medium
8	Kipper	200 cal		120 cal	Low calorie
9	NaN	NaN		NaN	NaN
10	Liver	200 cal		150 cal	Medium
11	Liver pate	150 cal		300 cal	Medium
12	Lamb (roast)	300 cal		300 cal	Med-High
13	Lobster boiled	200 cal		100 cal	Low calorie
14	NaN	NaN		NaN	NaN
15	Luncheon meat	300 cal		400 cal	High
16	Mackerel	320 cal		300 cal	Medium
17	Mussels	90 cal		90 cal	Low-Med
18	Pheasant roast	200 cal		200 cal	Medium
19	Pilchards (tinned)	140 cal		140 cal	Medium
20	Prawns	180 cal		100 cal	Low- Med
21	Pork	320 cal		290 cal	Med-High
22	Pork pie	320 cal		450 cal	High
23	Rabbit	200 cal		180 cal	Medium
24	Salmon fresh	220 cal		180 cal	Medium
25	Sardines tinned in oil	220 cal		220 cal	Medium
26	Sardines in tomato sauce	180 cal		180 cal	Medium
27	Sausage pork fried	250 cal		320 cal	High
28	Sausage pork grilled	220 cal		280 cal	Med-High
29	Sausage roll	290 cal		480 cal	High
30	Scampi fried in oil	400 cal		340 cal	High
31	Steak & kidney pie	400 cal		350 cal	High

```
# use list comprehension to convert the dataframe into a JSON string
tables_json = [table.to_json() for table in tables_df]
```

```
# Loop over each JSON string to print data from the table
for idx, table_json in enumerate (tables_json) :
    print(f"Table {idx + 1}:")
    print(table_json)
    # add a space/newline between tables
    print()
```

Table 1:

```
{"Fish cake":{"0":"Fish fingers","1":"Gammon","2":"Haddock fresh","3":
```

```
# extract tables from all pages
tables = read_pdf (pdf_file, pages='all', multiple_tables=True)
# print the tables extracted from each page
print(tables)
```

	BREADS & CEREALS	Portion size * \
0	Bagel (1 average)	140 cals (45g)
1	Biscuit digestives	86 cals (per biscuit)
2	Jaffa cake	48 cals (per biscuit)
3	Bread white (thick slice)	96 cals (1 slice 40g)
4	Bread wholemeal (thick)	88 cals (1 slice 40g)
5	Chapatis	250 cals
6	Cornflakes	130 cals (35g)
7	Crackerbread	17 cals per slice
8	Cream crackers	35 cals (per cracker)
9	Crumpets	93 cals (per crumpet)
10	Flapjacks basic fruit mix	320 cals
11	Macaroni (boiled)	238 cals (250g)
12	Muesli	195 cals (50g)
13	Naan bread (normal)	300 cals (small plate size)
14	Noodles (boiled)	175 cals (250g)
15	Pasta (normal boiled)	330 cals (300g)
16	Pasta (wholemeal boiled)	315 cals (300g)
17	Porridge oats (with water)	193 cals (350g)
18	Potatoes** (boiled)	210 cals (300g)
19	Potatoes** (roast)	420 cals (300g)

	per 100 grams (3.5 oz)	Unnamed: 0	energy content
0	310 cals	NaN	Medium
1	480 cals	NaN	High
2	370 cals	NaN	Med-High
3	240 cals	NaN	Medium

4	220 cal	NaN	Low-med	
5	300 cal	NaN	Medium	
6	370 cal	NaN	Med-High	
7	325 cal	NaN	Low Calorie	
8	440 cal	NaN	Low / portion	
9	198 cal	NaN	Low-Med	
10	500 cal	NaN	High	
11	95 cal	NaN	Low calorie	
12	390 cal	NaN	Med-high	
13	320 cal	NaN	Medium	
14	70 cal	NaN	Low calorie	
15	110 cal	NaN	Low calorie	
16	105 cal	NaN	Low calorie	
17	55 cal	NaN	Low calorie	
18	70 cal	NaN	Low calorie	
19	140 cal	NaN	Medium	
0	NaN	NaN	NaN	Rice (whit
1	Rice (egg-fried)	500 cal	200 cal	NaN
2	Rice (Brown)	405 cal (300g)	135 cal	NaN
3	Rice cakes	28 Cals = 1 slice	373 Cals	NaN
4	Ryvita Multi grain	37 Cals per slice	331 Cals	NaN
5	Ryvita + seed & Oats	180 Cals 4 slices	362 Cals	NaN
6	Spaghetti (boiled)	303 cal (300g)	101 cal	NaN
0	Low calorie			
1	High in portion			
2	Low calorie			
3	Medium			
4	Medium			
5	Medium			
6	Low calorie			
		Unnamed: 0	Unnamed: 1	

```
# set flag to process information page by page, performance optimizer
stream_option = True

# extract contents from page 4
page_number = 4

# extract tables in a rectangular area defined by coordinates (top, left,
area = (270, 13, 790, 900)

# extract from the specified area using the stream option
tables_df = read_pdf(pdf_file, pages=page_number, stream=stream_option, are

# loop over the table, print the information
for idx, table in enumerate (tables_df) :
    print (f"Table {idx + 1}:")
    print (table)
```

Table 1:

	Fruits & Vegetables	Portion size *	oz)	energy content
0	Apple	44 calories	44 calories	Low calorie
1	Banana	107 cal	65 calories	Low calorie
2	Beans baked beans	170 cal	80 calories	Low calorie
3	Beans dried (boiled)	180 cal	130 calories	Low calorie
4	Blackberries	25 cal	25 calories	Low calorie
5	Blackcurrant	30 cal	30 calories	Low calorie
6	Broccoli	27 cal	32 cal	Very low
7	Cabbage (boiled)	15 calories	20 calories	Low calorie
8	Carrot (boiled)	16 calories	25 calories	Low calorie
9	Cauliflower (boiled)	20 calories	30 calories	Low calorie
10	Celery (boiled)	5 calories	10 calories	Low calorie
11	Cherry	35 calories	50 calories	Low calorie
12	Courgette	8 cal	20 cal	Very low cal
13	Cucumber	3 calories	10 calories	Low calorie
14	Dates	100 calories	235 calories	Med-High
15	Grapes	55 calories	62 calories	Low calorie
16	Grapefruit	32 calories	32 calories	Low calorie
17	Kiwi	40 calories	50 calories	Low calorie
18	Leek (boiled)	10 calories	20 calories	Low calorie

Start coding or [generate](#) with AI.

Start coding or [generate](#) with AI.

