# DIGITAL ASSIGNMENT - 5

MAY 13, 2021
RISHABH SHARMA
20MAI0082

Github Link :- https://github.com/rishabh5197/Data-Mining/tree/main/Assignment-5

# Name :- Rishabh Sharma

# Registration Number :- 20MAI0082

# Assignment - 5

# Web Mining

### Importing Libraries

In [1]:
```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from bs4 import BeautifulSoup as bs
import requests
from warnings import filterwarnings
filterwarnings("ignore")
```

### Using Web Scraping to collect the data

In [2]:
```python
while True:
    try:
        link = 'https://www.flipkart.com/search?q=laptops&otracker=search&otracker1=se
        print(link)
        page = requests.get(link)
        print(page)
        break
    except:
        pass
```

https://www.flipkart.com/search?q=laptops&otracker=search&otracker1=search&marketpla
ce=FLIPKART&as-show=on&as=off (https://www.flipkart.com/search?q=laptops&otracker=s
earch&otracker1=search&marketplace=FLIPKART&as-show=on&as=off)
<Response [200]>

In [3]:
```python
soup = bs(page.content,'html.parser')
```

In [4]:
```python
# items = soup.find_all('div',class_="_4rR01T")
# items=[i.get_text() for i in items]
# items = [i.split(" - (")[0] for i in items]
```

In [5]:
```python
details = soup.find_all("li",class_="rgWa7D")
details = [i.get_text() for i in details]
```

```python
In [6]:   1  lis,newlis=[],[]
          2  for i in details:
          3      if ("Intel ")in i or ("AMD ") in i or ("M1") in i :
          4          if newlis:
          5              lis.append(newlis)
          6              newlis=[]
          7              newlis.append(i)
          8          else:
          9              newlis.append(i)
         10      else:
         11          newlis.append(i)
         12  newlisnew=[]
         13  for i in lis:
         14      newlisnew.append(" ".join(i))
```

```python
In [7]:   1  items = []
          2  for i in newlisnew:
          3      if ' GB DDR4' in i:
          4          items.append(i.split(" GB DDR4")[0][:-2])
          5      elif ' GB DDR3' in i:
          6          items.append(i.split(" GB DDR3")[0][:-2] )
          7      elif ' GB LPDDR4X' in i:
          8          items.append(i.split(" GB LPDDR4X")[0][:-2] )
          9      else:
         10          items.append("")
         11  company = [i.split()[0] for i in items]
```

```python
In [8]:   1  pages= soup.find_all('a',class_='ge-49M',href=True)
          2  pages = [str(i).split(">")[0][31:] for i in pages ]
```

```python
In [9]:   1  ratings = soup.find_all("div",class_='_3LWZlK')
          2  ratings = [i.get_text() for i in ratings]
```

```python
In [10]:  1  prices =  soup.find_all("div",class_="_30jeq3 _1_WHN1")
          2  prices = [i.get_text().replace("₹","").replace(",","") for i in prices]
```

```python
In [11]:  1  # details = soup.find_all("li",class_="rgWa7D")
          2  # details = [i.get_text() for i in details]
```

```python
In [12]:  1  # indepth_link = soup.find_all("a",class_ = '_1fQZEK')
```

```python
In [13]:  1  # start= str(indepth_link[0]).index("href=")
          2  # end = str(indepth_link[0]).index(" rel=")
          3  # links = ['https://www.flipkart.com'+str(str(i)[start+6:end-1]) for i in indepth_link]
```

In [14]:
```python
lis,newlis=[],[]
for i in details:
    if ("Intel ")in i or ("AMD ") in i or ("M1") in i :
        if newlis:
            lis.append(newlis)
            newlis=[]
            newlis.append(i)
        else:
            newlis.append(i)
    else:
        newlis.append(i)
newlisnew=[]
for i in lis:
    newlisnew.append(" ".join(i))
```

## Collecting data and filtering everything

In [15]:

```python
processor_brand=[]
processor = []
ram = []
ram_type=[]
os =[]
screen_size=[]
ssd_present=[]
ssd_capacity=[]
hdd_capacity=[]
# count=0
for i in newlisnew:
    if "Intel " in i:
#        print(count)
        processor_brand.append("Intel")
        processor.append(i.split()[1]+" "+i.split()[2])
    elif 'AMD' in i:
#        print(count)
        processor_brand.append("AMD")
        processor.append(i.split()[1]+" "+i.split()[2])
    elif "M1" in i:
#        print(count)
        processor_brand.append("M1")
        processor.append(i.split()[1]+" "+i.split()[2])
    else:
        processor_brand.append("")
        processor.append("")
#    count+=1
    if ' GB DDR4 RAM' in i:
        index = i.index(' GB DDR4 RAM')
        ram.append(i[index-2:index])
        ram_type.append("DDR4")
    elif ' GB DDR3' in i:
        index = i.index(' GB DDR3 RAM')
        ram.append(i[index-2:index])
        ram_type.append("DDR3")
    else:
        ram.append("")
    if ' Operating System' in i:
        a = i.split()
        index = a.index("Operating")
        if a[index-1]=="10":
            os.append("Windows 10")
        elif a[index-2]=='Mac':
            os.append("Mac Os")
        else:
            os.append(a[index-1])
    else:
        os.append("")
    if "inch" in i:
        a = i.split()
        if 'inches)' in a:
            index = a.index("inches)")
        else:
            index = a.index("inch)")
        screen_size.append(a[index-1].strip("(")+" inch")
    if (" GB SSD" in i) or (" TB SSD" in i):
```

```python
57          ssd_present.append("Yes")
58          if " GB SSD" in i:
59              start = i.index(" GB SSD")-3
60              end=i.index(" GB SSD")
61              ssd_capacity.append(i[start:end])
62          elif' TB SSD' in i:
63              index = i.index(" TB SSD")
64              ssd_capacity.append(int(i[index-1])*1024)
65      else:
66          ssd_present.append("No")
67          ssd_capacity.append("")
68      if " HDD" in i:
69          index = i.index(" HDD")
70          hdd_capacity.append(i[index-4])
71      else:
72          hdd_capacity.append("")
```

## Creating a dataframe in order to carry out further steps

```python
In [16]: 1  data = pd.DataFrame({
         2  'items':items[:len(processor)],                      # Name of Laptop
         3  'company':company[:len(processor)],                  # Company of laptop  (or Laptop
         4  'ratings out of 5':ratings[:len(processor)],          # What are the rating mentione
         5  'prices':prices[:len(processor)],                    # Price of the laptop
         6  'processor_brand':processor_brand,                   # what is the processor of the l
         7  'processor':processor,                               # Type of processor of the lapto
         8  'ram':ram ,                                          # how much ram does it have (Co
         9  'ram_type':ram_type,                                 # what is the type of ram (Categ
        10  'operating_system' :os,                              # consist of which operating sys
        11  'screen_size':screen_size,                           # Screen-size of the laptop (Ca
        12  'ssd_present':ssd_present,                           # is SSD present in the laptop (
        13  'ssd_capacity':ssd_capacity,                         # What is the capacity of the S
        14  'hdd_capacity_in_TB':hdd_capacity,                   # What is the capacity of HDD
        15  'Purchased': [np.random.choice(["No","Yes"]) for i in range(0,len(processor))]})  # this i
        16  data.to_csv("Dataset.csv",index=False)                                            # St
```

```python
In [17]: 1  data.shape
```

Out[17]: (23, 14)

```python
In [18]: 1  data =data.replace("",np.nan)
```

In [19]: 
```
1  data.head()
```

Out[19]:

| | items | company | ratings out of 5 | prices | processor_brand | processor | ram | ram_type | operating_s |
|---|---|---|---|---|---|---|---|---|---|
| 0 | AMD Ryzen 5 Quad Core Processor (3rd Gen) | AMD | 4.4 | 48990 | AMD | Ryzen 5 | 8 | DDR4 | Windo |
| 1 | Intel Core i3 Processor (10th Gen) | Intel | 4 | 35990 | Intel | Core i3 | 8 | DDR4 | Windo |
| 2 | Intel Core i5 Processor (9th Gen) | Intel | 4.5 | 52990 | Intel | Core i5 | 8 | DDR4 | Windo |
| 3 | Intel Core i3 Processor (10th Gen) | Intel | 4.2 | 33490 | Intel | Core i3 | 4 | DDR4 | Windo |
| 4 | Intel Core i3 Processor (10th Gen) | Intel | 4.2 | 35990 | Intel | Core i3 | 8 | DDR4 | Windo |

In [20]: 
```
1  data.isnull().sum()
```

Out[20]:
```
items                0
company              0
ratings out of 5     0
prices               0
processor_brand      0
processor            0
ram                  0
ram_type             0
operating_system     0
screen_size          0
ssd_present          0
ssd_capacity         4
hdd_capacity_in_TB   14
Purchased            0
dtype: int64
```

## Filling null values

In [21]:
```
1  data = data.replace(np.nan,0)
2  data.head()
```

Out[21]:

| | items | company | ratings out of 5 | prices | processor_brand | processor | ram | ram_type | operating_s |
|---|---|---|---|---|---|---|---|---|---|
| 0 | AMD Ryzen 5 Quad Core Processor (3rd Gen) | AMD | 4.4 | 48990 | AMD | Ryzen 5 | 8 | DDR4 | Windo |
| 1 | Intel Core i3 Processor (10th Gen) | Intel | 4 | 35990 | Intel | Core i3 | 8 | DDR4 | Windo |
| 2 | Intel Core i5 Processor (9th Gen) | Intel | 4.5 | 52990 | Intel | Core i5 | 8 | DDR4 | Windo |
| 3 | Intel Core i3 Processor (10th Gen) | Intel | 4.2 | 33490 | Intel | Core i3 | 4 | DDR4 | Windo |
| 4 | Intel Core i3 Processor (10th Gen) | Intel | 4.2 | 35990 | Intel | Core i3 | 8 | DDR4 | Windo |