# Insertion des ProductAttribute dans la table productAttribute

### Import des bibliothèques

Entrée [ ]:
```python
1  from lxml import etree
2  import pandas as pd
3  import mysql.connector
4  from mysql.connector import Error
5  import re
6  import json
```

### Parsage du catalogue Attribute

Entrée [ ]:
```python
1  # à modifier avant chaque traitement d'un nouveau fichier XML
2  refPath = 'unzipped_files/cat-attributes-fr-25-11-2019_11-41-09'
3
4  xtree = etree.parse(refPath)
5  xroot = xtree.getroot()
```

Entrée [ ]:
```python
1  for child in list(xroot):
2      print(child.tag)
```

### Génération d'un dataframe dfProdAtt

Entrée [ ]:
```python
1  df_cols = ["productId", "attributes"]
2  rows=[]
```

Entrée [ ]:
```python
1  for att in xroot.findall('products'):
2      #print(att.tag)
3      for i,att1 in enumerate(att.getchildren()):
4          #print(att1.attrib['productId'])
5          b = att1.getchildren()
6          for i,att2 in enumerate(b):
7              c = att2.getchildren()
8              row = []
9              for i,y in enumerate(c):
10                 # row = []
11                 row.append({"id": y.attrib['id'], "value": y.attrib['value']})
12             row_json = json.dumps(row, ensure_ascii=False)
13             rows.append({"productId": str(att1.attrib['productId']),
14                          "attributes": (row_json)})
15
16 dfProdAtt = pd.DataFrame(rows, columns=df_cols)
17 dfProdAtt.sort_values(['attributes'], ascending=False)
```

### NETTOYAGE DU DATAFRAME¶

#### Valeurs de productId dans XML cat

Entrée [ ]:
```python
1  # à modifier avant chaque traitement d'un nouveau fichier XML
2  refPath1 = 'unzipped_files/cat-ref-FR3787ED_2019-11-22 11.30.32-22-11-2019_11-48-2
3
4  xtree1 = etree.parse(refPath1)
5  xroot1 = xtree1.getroot()
```

Entrée [ ]:
```python
1  df_cols1 = ["productId"]
2  rows1 = []
3  for products in xroot1.iter('products'):
4      for a,ec in enumerate(products.getchildren()):
5          rows1.append(ec.attrib['productId'])
6  print(len(rows1))
```

### Recherche des doublons productId

Entrée [ ]:
```python
1  rows2 = []
2  doublons = []
3  for i in rows1:
4      if i not in rows2:
5          rows2.append(i)
6      else:
7          doublons.append(i)
8  print(len(rows2))
9  print(len(doublons))
```

### On supprime les doublons

Entrée [ ]:
```python
1  rows2 = list(set(rows2) - set(doublons))
2  print(len(rows2))
```

### Comparaison des valeurs de productId entre df et Product

Entrée [ ]:
```python
1  integrite = pd.Series(dfProdAtt['productId'].isin(rows2))
2  integrite
```

Entrée [ ]:
```python
1  #Insertion de la série integrité dans dfImage
2  dfProdAtt['Integrité'] = integrite
3  dfProdAtt.sort_values(['Integrité'], ascending=False)
```

Entrée [ ]:
```python
1  dfProdAtt['Integrité'].value_counts()
```

Entrée [ ]:
```python
1  resultat = dfProdAtt.groupby('productId').nunique()
2  resultat.shape
```

Entrée [ ]:
```python
1  dfProdAtt.attributes.replace("'"," ", regex=True, inplace=True)
2  dfProdAtt.attributes.replace("\\\\","", regex=True, inplace=True)
3  dfProdAtt.attributes.replace('1/2'","1/2", regex=True, inplace=True)
4  dfProdAtt.attributes.replace('3/4'","3/4", regex=True, inplace=True)
5  dfProdAtt.attributes.replace('3/8'","3/8", regex=True, inplace=True)
6  dfProdAtt.attributes.replace('1/4'","1/4", regex=True, inplace=True)
```

```
1  dfProdAtt.attributes.replace('1/2}','1/2"}', regex=True, inplace=True)
2  dfProdAtt.attributes.replace('3/4}','3/4"}', regex=True, inplace=True)
3  dfProdAtt.attributes.replace('3/8}','3/8"}', regex=True, inplace=True)
4  dfProdAtt.attributes.replace('1/4}','1/4"}', regex=True, inplace=True)
5  dfProdAtt.attributes.replace('18.90"','18.90', regex=True, inplace=True)
6  dfProdAtt.attributes.replace('17.72"','17.72', regex=True, inplace=True)
```

## Insertion des données du dataframe dans la table Attribute

```
1  connection_config = {
2                       'host':"localhost",
3                       'port': 3308,
4                       'database': 'bihr_db',
5                       'user': 'BASTIER',
6                       'passwd': "DA2019",
7                       #'autocommit': True
8                       }
```

```
1  try:
2      connection = mysql.connector.connect(**connection_config)
3
4      for i in range(dfProdAtt.shape[0]):
5          prodAtt = dfProdAtt.iloc[i]
6          print(i)
7          if prodAtt['Integrité'] == True :
8              print(prodAtt['Integrité'])
9              print(prodAtt['attributes'])
10             ProdAttInsertQuery = """INSERT INTO productattribute (productId, attri
11                             VALUES
12                             ("""+ """'""" + str(prodAtt['productId']) + """','"""+
13             cursor = connection.cursor()
14             result = cursor.execute(ProdAttInsertQuery)
15         #print(prodAtt['productId'])
16
17     connection.commit()
18     print("Insertion datas in ProdAtt table successful ")
19     cursor.close()
20
21 except mysql.connector.Error as error:
22     print("Failed to insert datas in ProdAtt table : {}".format(error))
23
24 finally:
25     if (connection.is_connected()):
26         cursor.close()
27         connection.close()
28         print("MySQL connection is closed")
```

```
1
```