

TP sur les files et les piles

Thème 1

I Motivation

1 Fichier au format HTML bien formés

Le HTML est un format de fichier utilisé par les navigateurs web. Les fichiers au format HTML (et plus généralement au format XML) sont des fichiers texte dans lesquels on trouve des balises :

- ouvrantes de la forme `<nom attributs>`,
- fermantes de la forme `</nom>`,
- auto-fermantes de la forme `<nom attributs/>`.

où `nom` désigne le nom de la balise et `attributs` une liste de couples `clé=valeur`.

Dans la pratique, `nom` est par exemple `div`, `p`, `html`, `body`, `head`, ... Pour la syntaxe des tags HTML, voir la page du W3C dédiée¹. Les balises auto-fermantes sont : `area`, `br`, `hr`, `img`, `input`, `link`, `meta` et `param`.

Dans cette activité, on considère qu'un document HTML est bien formé si :

- à chaque balise ouvrante correspond une balise fermante,
- on ne peut fermer une balise que si toutes les balises situées entre les deux balises ouvrantes et fermantes sont fermées.

Par exemple, les documents `ex1.html` et `ex4.html` sont bien formés, alors que les documents `ex2.html` et `ex3.html` sont mal formés.

On veut écrire un prédicat² renvoyant `True` si un texte est un document HTML bien formé et `False` dans le cas contraire.

Dans ce cadre, on peut donc ne pas tenir compte des balises auto-fermantes (et on ne demande pas de vérifier ici si les balises sont des balises HTML existantes, ni si on respecte les attributs des balises).

La première étape est d'écrire un parser³ de fichier HTML, permettant de parcourir séquentiellement les balises.

II Écriture d'un parser HTML

L'écriture d'un tel parser est une tâche difficile, car un caractère '`<`' peut être rencontré dans différentes situations :

- Situation 1 : définition du type du document : `<!DOCTYPE ...>`
- Situation 2 : signe inférieur dans le texte : `i < len(1)`
- Situation 3 : commentaires HTML `<!--` et `-->`
- Situation 4 : signe inférieur dans des attributs : `<script data-user=">myfunc();">`

Pour simplifier les choses nous commencerons par traiter uniquement des fichiers HTML avec la situation 1, puis nous autoriserons dans une deuxième version les documents HTML contenant les situations 2 à 4.

1 Structure de file

Les tags rencontrés seront ajoutés dans une file. Une file est une structure de données séquentielle agissant comme une file d'attente : on peut ajouter et retirer des éléments de la structure, mais les premiers éléments ajoutés seront toujours les premiers à sortir.

La structure de file dispose des primitives suivantes :

- `q=Queue()` : crée une nouvelle file vide;
- `q.enqueue(e1)` : ajoute l'élément `e1` à la file `q`;
- `q.dequeue()` : enlève le premier élément ajouté à la structure et le renvoie;
- `q.is_empty()` : renvoie `True` si, et seulement si, la file est vide.

1. https://www.w3.org/community/webed/wiki/HTML/Training/Tag_syntax

2. **Prédicat** : dans ce contexte, il s'agit d'une fonction retournant `True` ou `False` en relation à une propriété du paramètre de cette dernière.

3. **parser** : Programme informatique qui met en évidence la structure d'un texte, d'un programme informatique.

On dit que la file est une structure **FIFO** (*First In First Out*).

Le fichier `myqueue.py` contient le squelette d'un module implémentant une structure de file (on ne peut utiliser le nom `queue.py` car Python dispose déjà d'un module homonyme).

A faire

1. Prendre connaissance du paragraphe 5.1.2 de la documentation Python ^a.
Comme ce paragraphe l'indique : *toutefois, les listes ne sont pas très efficaces pour réaliser ce type de traitement. Alors que les ajouts et suppressions en fin de liste sont rapides, les opérations d'insertions ou de retraits en début de liste sont lentes (car tous les autres éléments doivent être décalés d'une position).*
Aussi, il sera judicieux d'utiliser comme attribut de notre classe `Queue` une structure de la classe `collections.deque`.
2. Recopier le fichier `myqueue.py` dans votre espace de travail et compléter-le.

a. 5.1.2 Utilisation des listes comme des files : <https://docs.python.org/fr/3/tutorial/datastructures.html#using-lists-as-queues>

2 Première version

Dans un module `html_parser1.py`, écrire une première version du parser : la fonction `parse(document)` prend en paramètre un document sous forme d'une chaîne de caractères et renvoie une file contenant les tags.

Pour vous aider, le fichier `html_parser1.py` contient le squelette d'un parser.

Dans cette première version, on considère que seule la situation 1 peut être rencontrée et on ignore les situations 2, 3 et 4. Et on rappelle que les balises auto-fermantes peuvent être ignorées.

Par souci de simplification on pourra considérer que les attributs d'une balise ne peuvent être séparés que par des espaces.

On utilisera à profit la méthode `index` ⁴ des chaînes de caractères :

```

1  >>> help(str.index)
2  Help on method_descriptor:
3
4  index(...)
5      S.index(sub[, start[, end]]) -> int
6
7      Return the lowest index in S where substring sub is found,
8      such that sub is contained within S[start:end]. Optional
9      arguments start and end are interpreted as in slice notation.
10
11     Raises ValueError when the substring is not found.
```

La méthode `split` ⁵ peut également s'avérer utile.

A faire

1. Recopier le fichier `html_parser1.py` dans votre espace de travail et compléter-le en tenant des conseils ci-dessus.

III Vérificateur HTML : un algorithme

Dès que l'on peut récupérer les tags séquentiellement, nous pouvons nous attaquer à l'écriture du vérificateur HTML (checker). Il peut être utile de faire l'analogie entre les documents HTML et les expressions correctement parenthésées.

Une expression bien parenthésée est une expression contenant un certains nombre de type de parenthèses telle que :

- à chaque parenthèse ouvrante correspond une parenthèse fermante,
- à tout moment, on ne peut fermer une parenthèse que si l'expression située entre les deux parenthèses se correspondant est bien parenthésée.

4. Méthode `index` : <https://docs.python.org/fr/3/library/stdtypes.html?highlight=index#str.index>

5. Méthode `split` : <https://docs.python.org/3.6/library/stdtypes.html#str.split>

Par exemple, les expressions suivantes sont bien parenthésées :

- ((()))
- ()()()
- ([{}])()

Les expressions suivantes ne sont pas bien parenthésées :

- (())
- (()
- ([)]

Pour vérifier si une expression est correctement parenthésée, une méthode est d'utiliser une pile :

- lorsque l'on rencontre une parenthèse ouvrante on l'empile,
- lorsque l'on rencontre une parenthèse fermante, par exemple] :
 - si la pile est vide, alors l'expression est mal parenthésée (trop de fermantes),
 - sinon, on dépile l'ouvrante située au sommet de la pile et on vérifie que les deux parenthèses correspondent ((et), [et]).

Lorsque toute la chaîne a été parcourue, la pile doit être vide (pas d'ouvrante non fermée).

Un document HTML peut être vue comme une expression parenthésée :

- les tags ouvrants '<tag>' sont les parenthèses ouvrantes;
- les tags fermants '</tag>' sont les parenthèses fermantes.

Sur la pile, on placera des tags ouvrants.

IV Opération primitives sur les piles

Les *piles* sont des structures de données linéaires admettant les *opérations primitives* suivantes :

- création d'une pile vide;
- empilement d'un élément sur une pile;
- dépilement du sommet d'une pile;
- test de vacuité d'une pile.

Le principe d'une pile est le suivant : on ne peut accéder un élément (e) de la pile qu'en ayant enlevé d'abord tous les éléments empilés après (e).

L'élément renvoyé par la méthode pop est le dernier élément empilé. On dit que la pile est une structure **LIFO** (Last In First Out).

V Implémentation du module mystack

Le fichier `mystack.py` contient le squelette de l'implémentation d'une pile.

A faire

1. Copier dans votre espace de travail le fichier `mystack.py` et implémentez votre propre structure de pile.
2. Enfin, copier dans votre espace de travail le fichier `html_checker.py` et implémenter l'algorithme de vérification.

VI Amélioration du parser

1 Utiliser des expressions régulières

L'objectif est ici d'écrire une nouvelle version du parser permettant de prendre en compte l'ensemble des situations particulières listées plus haut.

Pour cela, nous pouvons utiliser une expression régulière. Une expression régulière est une chaîne de caractère spécifiant un format permettant de :

- dire si une chaîne correspond au format;
- capturer certaines parties de ce format.

Un cours sur les expressions régulières est hors du programme, mais nous allons en fournir une pour récupérer les tags d'un document html.

Le code du fichier `html_parser2.py` contient une expression régulière présentée sur plusieurs lignes pour plus de lisibilité. Elle permet de récupérer les tags en tenant compte des guillemets.

A faire

- Copier le fichier `html_parser2.py` dans votre espace de travail.
- Bien regarder l'exemple d'utilisation de l'expression régulière pour la comprendre (lignes 30 à 45).
Ne pas hésiter à tester plusieurs exemples dans la console Python.
- En utilisant cette expression régulière, écrire une deuxième version du parser en complétant la fonction `parser`.

2 Pour aller plus loin

Utiliser une expression régulière a toujours ses limites pour parser un fichier html. Des bibliothèques existent pour itérer sur les tags d'un tel document.

Parmi elles, on peut utiliser la classe `HTMLParser`. Le fichier `html_parser3.py` contient le code d'une classe héritant de `HTMLParser` et qui construit la liste des tags.

A faire

1. Copier le fichier `html_parser3.py` dans votre espace de travail.
2. Compléter la fonction `parser` utilisant la classe décrite dans ce même fichier.

VII Source

Université de Lille (2020, 21 novembre). *Piles et files*. Enseignement DIU EIL.

https://gitlab-fil.univ-lille.fr/diu-eil-lil/portail/-/blob/master/bloc4-5/pile_file/pile.md

VIII Annexes

1 myqueue.py

```
1  import collections
2
3  class QueueEmptyError(Exception):
4      def __init__(msg):
5          super().__init__(msg)
6
7  class Queue:
8      """
9      a class representing a queue
10
11      >>> q = Queue()
12      >>> q.is_empty()
13      True
14      >>> q.enqueue(1)
15      >>> q.enqueue(2)
16      >>> q.is_empty()
17      False
18      >>> q.dequeue()
19      1
20      >>> q.dequeue()
21      2
22      """
23
24      def __init__(self):
25          """
26          create a new queue
27          """
28          self.inside = #A compléter avec un objet vide de la classe "deque"
29
30      def is_empty(self):
31          """
32          :return: (bool) True si la queue est vide, False sinon
33          """
34          pass
35
36      def enqueue(self, el):
37          """
38          enfile un élément dans la file
39          :param el: (any) un élément
40          """
41          pass
42
43      def dequeue(self):
44          """
45          défile un élément
46          :return: (any) un élément
47          """
48          pass
49
50  if __name__ == "__main__":
51      import doctest
52      doctest.testmod(verbose = True)
```

2 html_parser1.py

```
1  from myqueue import Queue
2
3  def parse(doc):
4      """
5          :param doc: (str) a document to be parsed
6          :return: (Queue) a queue containing tag (without attributes)
7          :CU: None
8          :Example:
9
10         >>> q = parse("<!DOCTYPE html><html><div class='titre'>du contenu ignoré</div></html>")
11         >>> q.is_empty()
12         False
13         >>> q.dequeue() == "<html>"
14         True
15         >>> q.dequeue() == "<div>"
16         True
17         >>> q.dequeue() == "</div>"
18         True
19         >>> q.dequeue() == "</html>"
20         True
21         >>> q.is_empty()
22         True
23         >>> # Other example
24         >>> q = parse("<html><div><p <a")
25         >>> q.dequeue() == "<html>"
26         True
27         >>> q.dequeue() == "<div>"
28         True
29         >>> q.is_empty()
30         True
31         """
32     pass
33
34 if __name__ == "__main__":
35     import doctest
36     doctest.testmod()
```

3 mystack.py

```
1 class StackEmptyError(Exception):
2     """
3     exception pour pile vide
4     """
5     def __init__(self, msg):
6         self.message = msg
7
8 class Stack:
9     """
10    une classe pour manipuler les piles
11    """
12
13    def __init__(self):
14        """
15        constructeur de pile
16        """
17        pass
18
19    def is_empty(self):
20        """
21        :return: (bool) True si la pile est vide, False sinon
22        :CU: None
23        :Exemples:
24
25        >>> p = Stack()
26        >>> p.is_empty()
27        True
28        >>> p.push(1)
29        >>> p.is_empty()
30        False
31        """
32        pass
33
34    def push(self, el):
35        """
36        :param el: (any) un élément
37        :return: None
38        :Side-Effet: ajoute un élément au sommet de la pile
39        :CU: None
40
41        >>> p = Stack()
42        >>> p.push(1)
43        >>> p.pop() == 1
44        True
45        """
46        pass
```

```
1     def pop(self):
2         """
3         :return: (any) l'élément au sommet de la pile
4         :CU: la pile ne doit pas être vide
5         :raise: StackEmptyError
6         :Side-Effect: la pile est modifiée
7         :Exemples:
8
9         >>> p = Stack()
10        >>> p.push(1)
11        >>> p.pop() == 1
12        True
13        >>> p.is_empty()
14        True
15        """
16        pass
17
18    def top(self):
19        """
20        :return: (any) l'élément au sommet de la pile
21        :CU: la pile ne doit pas être vide
22        :raise: StackEmptyError
23        :Exemples:
24
25        >>> p = Stack()
26        >>> p.push(1)
27        >>> p.top() == 1
28        True
29        >>> p.is_empty()
30        False
31        """
32        pass
33
34    if __name__ == "__main__":
35        import doctest
36        doctest.testmod(verbose=True)
```

4 html_checker.py

```
1  from mystack import Stack
2  from myqueue import Queue
3  from html_parser1 import parse
4
5  def html_checker(s):
6      """
7      :param s: (str) un document html
8      :return: (bool) True si `s` est bien formé, False sinon
9      :CU: Aucune
10     :Exemple:
11
12     >>> html_checker("<!DOCTYPE html><html lang='fr'><div><p></p></div></html>")
13     True
14     >>> html_checker("<!DOCTYPE html><html lang='fr'><div></p></div><p></html>")
15     False
16     >>> html_checker("<!DOCTYPE html><html lang='fr'><div><p></p></div>")
17     False
18     """
19     pass
20
21  if __name__ == "__main__":
22     import sys
23     if len(sys.argv) <= 1:
24         # pas d'arguments sur la ligne de commande
25         # on exécute les tests.
26         import doctest
27         doctest.testmod(verbose = True)
28     else:
29         # on considère que l'argument est le nom d'un fichier
30         fname = sys.argv[1]
31         try:
32             with open(fname, 'r') as f:
33                 if html_checker(f.read()):
34                     print("le fichier {} est un fichier html correct.".format(fname))
35                 else:
36                     print("le fichier {} n'est pas un fichier html correct.".format(fname))
37         except:
38             print("Impossible de parser le fichier {}".format(fname))
```

5 html_parser2.py

```

1  from myqueue import Queue
2
3  import re
4
5  #Expression régulière
6  HTML_REGEX = re.compile(
7      r"""<                                # commence par un <
8      (/?\w+)                             # capture le nom du tag (éventuellement avec un slash
9      (                                     # groupe des attributs
10         (
11             \s+                           # autant d'espace que l'on veut
12             \w+                           # suivi d'un mot (nom de l'attribut)
13             (
14                 \s*=\s*                   # un signe égal, entouré d'espaces
15                 (?:".*?"|'.*?'|/\^[^"]>\s]+) # valeur de l'attribut : capture de tout ce qui est entre guillemets
16             )?                             # valeur optionnelle
17         )+                               # autant de couple nom=valeur que l'on veut
18         \s*                               # suivi d'éventuelles espaces
19         /\s*                             # ou bien pas d'attributs : que des espaces
20     )                                     # fin des attributs
21     >                                    # finit par un >
22 """, re.VERBOSE)
23
24 #Version courte :
25 '''
26 HTML_REGEX = re.compile(
27     r"""<(/?\w+)((\s+\w+(\s*=\s*(?:".*?"|'.*?'|/\^[^"]>\s]+))?)\s*/\s*>""")
28 '''
29
30 #Exemple d'utilisation :
31 """
32 >>> html = '<!DOCTYPE html> <!-- un commentaire --> <html lang="fr-FR"> <body>' +\
33         "<h1> Un titre </h1> <p class='code'> un paragraphe   while i < len(l):" +\
34         ' </p> <br/> </body> </html>'
35 >>> for tag_element in HTML_REGEX.findall(html):
36     print(tag_element[0])
37
38 html
39 body
40 h1
41 /h1
42 p
43 /p
44 /body
45 /html
46 """

```

```
1 def parse(doc):
2     """
3     :param doc: (str) a document to be parsed
4     :return: (Queue) a queue containing tag (without attributes)
5     :CU: None
6     :Example:
7
8     >>> q = parse("<!DOCTYPE html><html><div class='titre'>du contenu ignoré</div></html>")
9     >>> q.is_empty()
10    False
11    >>> q.dequeue() == "html"
12    True
13    >>> q.dequeue() == "div"
14    True
15    >>> q.dequeue() == "/div"
16    True
17    >>> q.dequeue() == "/html"
18    True
19    >>> q.is_empty()
20    True
21    >>> # Other example
22    >>> q = parse("<html><div><p <a")
23    >>> q.dequeue() == "html"
24    True
25    >>> q.dequeue() == "div"
26    True
27    >>> q.is_empty()
28    True
29    """
30    pass
31
32 if __name__ == "__main__":
33     import doctest
34     doctest.testmod()
```

6 html_parser3.py

```
1  from html.parser import HTMLParser
2
3  class MyHTMLParser(HTMLParser):
4      """
5          a class that parse a document and allow tag access
6          :examples:
7
8          >>> parser = MyHTMLParser("<!DOCTYPE hml><html lang='fr'></html>")
9          >>> parser.has_tag()
10         True
11         >>> parser.next_tag()
12         '<html>'
13         >>> parser.next_tag()
14         '</html>'
15         >>> parser.has_tag()
16         False
17         """
18     def __init__(self, data):
19         """
20             constructor for MyHTMLParse
21             :param data: (str) html document
22             :UC: None
23             """
24         super().__init__()
25         self.__tags = []
26         self.__tag_index = 0
27         HTMLParser.feed(self, data)
28     def handle_starttag(self, tag, attrs):
29         """
30             handle an opening tag
31             :param tag: (str) the opening tag
32             :param attrs: (list) attributes
33             """
34         self.__tags.append('<{:s}>'.format(tag))
35     def handle_endtag(self, tag):
36         """
37             handle an ending tag
38             :param tag: (str) the ending tag
39             """
40         self.__tags.append('</{:s}>'.format(tag))
41     def has_tag(self):
42         """
43             :return: (bool) True if document contains another tag, False otherwise
44             """
45         return self.__tag_index < len(self.__tags)
46     def next_tag(self):
47         """
48             :return: (str) the next tag in document
49             """
50         res = self.__tags[self.__tag_index]
51         self.__tag_index += 1
52         return res
53
54     def parser(doc):
55         pass
56
57     if __name__ == '__main__':
58         import doctest
59         doctest.testmod(optionflags=doctest.NORMALIZE_WHITESPACE | doctest.ELLIPSIS, verbose=False)
```
