

Федеральное государственное автономное образовательное учреждение высшего
образования
«Национальный исследовательский университет ИТМО»
Факультет программной инженерии и компьютерной техники

Отчет
по лабораторной работе №4
«Исследование протоколов, форматов обмена информацией
и языков разметки документов»
по дисциплине «Информатика»

вариант 35

Выполнил: Хоробрых Д.Е., группа Р3116
Преподаватель: Машина Е.А.

Санкт-Петербург
~ 2022 ~

Оглавление

Задание	3
Задание 1	4
Задание 2	8
Задание 3	9
Задание 4	10
Задание 5	12
Вывод.....	14
Список литературы	15

Задание

Вариант:

№ варианта	Исходный формат	Результирующий формат	День недели
35	XML	YAML	Суббота

- 1) Написать программу на языке Python 3.x, которая бы осуществляла парсинг и конвертацию исходного файла в новый. Готовые библиотеки и регулярные выражения использовать нельзя.
- 2) Найти готовые библиотеки, осуществляющие аналогичный парсинг и конвертацию файлов. Переписать исходный код, применив найденные библиотеки. Регулярные выражения использовать нельзя. Сравнить полученные результаты и объяснить их сходство/различие.
- 3) Переписать исходный код, добавив в него использование регулярных выражений. Сравнить полученные результаты и объяснить их сходство/различие.
- 4) Используя свои программы из заданий 1–3, сравнить стократное время выполнения парсинга + конвертации в цикле. Проанализировать полученные результаты и объяснить их сходство/различие.
- 5) Переписать исходную программу, чтобы она осуществляла парсинг и конвертацию исходного файла в любой другой формат (кроме JSON, YAML, XML, HTML): PROTOBUF, TSV, CSV, WML и т. п..

Задание 1

Код программы:

```
1. def parse_xml_to_yaml(filename: str):
2.     result = ''
3.     TAB = ' '
4.     with open(filename, 'r', encoding='utf8') as f:
5.         text = f.read()
6.         now_tab_level = 0
7.         for line in map(str.strip, text.split('\n')):
8.             writing_tag = False
9.             now_tag = ''
10.            last_symbol = ''
11.            string = ''
12.            if line.startswith('<?'):
13.                continue
14.            if line.startswith('</'):
15.                now_tab_level -= 1
16.                continue
17.            for index, symbol in enumerate(line):
18.                if symbol == '<':
19.                    writing_tag = True
20.                elif last_symbol == '<' and symbol == '/':
21.                    result += TAB * now_tab_level + now_tag + ': ' + string
22.                    break
23.                elif symbol == '>':
24.                    writing_tag = False
25.                    if index == len(line) - 1:
26.                        result += TAB * now_tab_level + now_tag + ':'
27.                        now_tab_level += 1
28.                elif writing_tag:
29.                    now_tag += symbol
30.                else:
31.                    string += symbol
32.                    last_symbol = symbol
33.            result += '\n'
34.        return result
35.
36.
37. with open('result.yaml', 'w', encoding='utf8') as f:
38.     f.write(parse_xml_to_yaml('day.xml'))
```

Содержание исходного файла day.xml:

```
1. <?xml version="1.0" encoding="UTF-8" ?>
2. <schedule>
3.   <day>
4.     <day_date>сб, 29 октября</day_date>
5.     <lesson1>
6.       <time>
7.         <start>11:40</start>
8.         <end>13:10</end>
9.       </time>
10.      <body>
11.        <type>Лекция</type>
12.        <name>Математика (базовый уровень)</name>
13.        <info>
14.          <teacher>Правдин Константин Владимирович</teacher>
15.          <place>Ауд. 1404, Кронверкский пр., д.49, лит.А</place>
16.          <distant>Очный</distant>
17.        </info>
18.      </body>
19.    </lesson1>
20.    <lesson2>
21.      <time>
22.        <start>11:40</start>
23.        <end>13:10</end>
24.      </time>
25.      <body>
26.        <type>Лекция</type>
27.        <name>Математика (базовый уровень)</name>
28.        <info>
29.          <teacher>Правдин Константин Владимирович</teacher>
30.          <place>Ауд. 2201 (бывш. 206), Кронверкский пр., д.49,
лит.А</place>
31.          <distant>Очно-дистанционный</distant>
32.        </info>
33.      </body>
34.    </lesson2>
35.    <lesson3>
36.      <time>
37.        <start>13:30</start>
38.        <end>15:00</end>
39.      </time>
40.      <body>
41.        <type>Лекция</type>
42.        <name>Математика (базовый уровень)</name>
43.        <info>
44.          <teacher>Правдин Константин Владимирович</teacher>
45.          <place>Ауд. 1404, Кронверкский пр., д.49, лит.А</place>
46.          <distant>Очный</distant>
```

```

47.         </info>
48.     </body>
49. </lesson3>
50. <lesson4>
51.     <time>
52.         <start>13:30</start>
53.         <end>15:00</end>
54.     </time>
55.     <body>
56.         <type>Лекция</type>
57.         <name>Математика (базовый уровень)</name>
58.         <info>
59.             <teacher>Правдин Константин Владимирович</teacher>
60.             <place>Ауд. 2201 (бывш. 206), Кронверкский пр., д.49,
лит.А</place>
61.             <distant>Очно-дистанционный</distant>
62.         </info>
63.     </body>
64. </lesson4>
65. <lesson5>
66.     <time>
67.         <start>11:40</start>
68.         <end>13:10</end>
69.     </time>
70.     <body>
71.         <type>Лекция</type>
72.         <name>Математика (базовый уровень)</name>
73.         <info>
74.             <teacher>Правдин Константин Владимирович</teacher>
75.             <place>Ауд. 1404, Кронверкский пр., д.49, лит.А</place>
76.             <distant>Очный</distant>
77.         </info>
78.     </body>
79. </lesson5>
80. </day>
81.</schedule>

```

Содержание результирующего файла result.yaml:

```

1. schedule:
2.   day:
3.     day_date: сб, 29 октября
4.     lesson1:
5.       time:
6.         start: 11:40
7.         end: 13:10
8.       body:
9.         type: Лекция

```

```

10.     name: Математика (базовый уровень)
11.     info:
12.         teacher: Правдин Константин Владимирович
13.         place: Ауд. 1404, Кронверкский пр., д.49, лит.А
14.         distant: Очный
15. lesson2:
16.     time:
17.         start: 11:40
18.         end: 13:10
19.     body:
20.         type: Лекция
21.         name: Математика (базовый уровень)
22.         info:
23.             teacher: Правдин Константин Владимирович
24.             place: Ауд. 2201 (бывш. 206), Кронверкский пр., д.49, лит.А
25.             distant: Очно-дистанционный
26. lesson3:
27.     time:
28.         start: 13:30
29.         end: 15:00
30.     body:
31.         type: Лекция
32.         name: Математика (базовый уровень)
33.         info:
34.             teacher: Правдин Константин Владимирович
35.             place: Ауд. 1404, Кронверкский пр., д.49, лит.А
36.             distant: Очный
37. lesson4:
38.     time:
39.         start: 13:30
40.         end: 15:00
41.     body:
42.         type: Лекция
43.         name: Математика (базовый уровень)
44.         info:
45.             teacher: Правдин Константин Владимирович
46.             place: Ауд. 2201 (бывш. 206), Кронверкский пр., д.49, лит.А
47.             distant: Очно-дистанционный
48. lesson5:
49.     time:
50.         start: 11:40
51.         end: 13:10
52.     body:
53.         type: Лекция
54.         name: Математика (базовый уровень)
55.         info:
56.             teacher: Правдин Константин Владимирович
57.             place: Ауд. 1404, Кронверкский пр., д.49, лит.А
58.             distant: Очный

```

Задание 2

Для выполнения данного задания были выбраны следующие готовые библиотеки:

Xml_to_dict - для парсинга xml файла в словарь Python.

PuYAML - для конвертации словаря Python в формат YAML.

Исходный код программы:

```
1. from xml_to_dict import XMLtoDict
2. from yaml import dump
3.
4.
5. def task2_with_libs(filename: str):
6.     with open(filename) as f:
7.         dictionary = XMLtoDict().parse(f.read())
8.     with open('result.yaml', 'w') as f:
9.         dump(dictionary, f, allow_unicode=True)
10.
11.
12.task2_with_libs('day.xml')
```

Исходный и результирующий файлы аналогичны заданию 1.

По сравнению с первым вариантом выполнения задания, использование готовых сторонних библиотек позволило использовать намного меньше строк кода для реализации поставленной задачи, а также сократило время на создание и кода.

Задание 3

Для выполнения данного задания были произведены некоторые манипуляции с исходным файлом с помощью регулярных выражений.

Исходный код программы:

```
1. import re
2.
3.
4. def parser_xml_to_yaml_re(filename: str):
5.     with open(filename, encoding='utf8') as f:
6.         res = re.sub(re.compile(r'<?\.*>'), '', f.read()) # удаление
# служебных тегов
7.         res = re.sub(re.compile(r'</\w+>'), '', res) # удаление
# закрывающихся тегов
8.         res = re.sub(re.compile(r'>'), r': ', res) # приведение тегов к
# виду YAML
9.         res = re.sub(re.compile(r'<'), '', res)
10.        res = '\n'.join(filter(lambda x: x.strip(), res.split('\n'))) #
# удаление пустых строк
11.
12.    with open('result.yaml', 'w', encoding='utf-8') as f:
13.        f.write(res)
14.
15.
16. parser_xml_to_yaml_re('day.xml')
```

Исходный и результирующий файлы аналогичны заданию 1.

По сравнению с первым вариантом выполнения задания для реализации задачи с помощью регулярных выражений потребовалось также гораздо меньше строк кода, однако понадобилось время для продумывания оптимальных шаблонов и их написания.

Задание 4

В данном задании необходимо сравнить время выполнения трех вариантов программ при парсинге и конвертации 100 файлов в цикле.

Исходный код программы:

```
1. import datetime
2.
3. from task1F import parse_xml_to_yaml
4. from task2 import task2_with_libs
5. from task3 import parser_xml_to_yaml_re
6.
7. results = list()
8.
9. start = datetime.datetime.now()
10. for i in range(100):
11.     parse_xml_to_yaml('day.xml')
12. results.append(datetime.datetime.now() - start)
13.
14. start = datetime.datetime.now()
15. for i in range(100):
16.     task2_with_libs('day.xml')
17. results.append(datetime.datetime.now() - start)
18.
19. start = datetime.datetime.now()
20. for i in range(100):
21.     parser_xml_to_yaml_re('day.xml')
22. results.append(datetime.datetime.now() - start)
23.
24. print('Собственный код: ' + str(results[0].total_seconds()))
25. print('Использование библиотек: ' + str(results[1].total_seconds()))
26. print('Использование регулярных выражений: ' +
        str(results[2].total_seconds()))
```

Исходный и результирующий файлы аналогичны заданию 1.

Вывод программы:

```
Собственный код: 0.028428
Использование библиотек: 0.44911
Использование регулярных выражений: 0.048081
```

Итак, было получено, что наиболее быстрым вариантом парсинга и конвертации большого количества файлов является вариант программы, написанный без сторонних библиотек и регулярных выражений. Следующим по затратам времени оказался вариант с регулярными выражениями, а самым медленным (почти в 10 раз медленнее остальных) - вариант с использованием сторонних библиотек.

Данный результат являлся вполне ожидаемым, потому что внешние библиотеки не оптимизированы между собой и для реализации поставленной задачи необходимо переводить исходный файл в формат словаря Python.

Самым быстрым оказался собственный парсер, потому что в нем не используются никакие сторонние модули, только стандартные методы строк.

Задание 5

В данном задании необходимо написать программу для конвертации XML файла в любой другой (кроме JSON, YAML, XML, HTML) формат.

Мною был выбран формат CSV.

Исходный код программы:

```
1. import csv
2.
3. from xml_to_dict import XMLtoDict
4.
5.
6. def task5(filename: str):
7.     with open(filename, 'r', encoding='utf8') as f:
8.         dictionary = XMLtoDict().parse(f.read())
9.         result = list()
10.        for key, value in dictionary.get('schedule').items():
11.            date = value['day_date']
12.            for k, v in value.items():
13.                if k == 'day_date':
14.                    date = v
15.                    continue
16.            start = v['time']['start']
17.            end = v['time']['end']
18.            type = v['body']['type']
19.            name = v['body']['name']
20.            teacher = v['body']['info']['teacher']
21.            place = v['body']['info']['place']
22.            distant = v['body']['info']['distant']
23.
24.            result.append({'date': date,
25.                           'start': start,
26.                           'end': end,
27.                           'type': type,
28.                           'name': name,
29.                           'teacher': teacher,
30.                           'place': place,
31.                           'distant': distant})
32.
33.        with open('result.csv', 'w', encoding='utf8') as f:
34.            d_writer = csv.DictWriter(f, result[0].keys())
35.            d_writer.writeheader()
36.            d_writer.writerows(result)
37.
38.
39. task5('day.xml')
```

Исходный файл аналогичен заданию 1.

Результирующий файл result.csv:

```
date,start,end,type,name,teacher,place,distant
"сб, 29 октября",11:40,13:10,Лекция,Математика (базовый уровень),Правдин
Константин Владимирович,"Ауд. 1404, Кронверкский пр., д.49, лит.А",Очный
"сб, 29 октября",11:40,13:10,Лекция,Математика (базовый уровень),Правдин
Константин Владимирович,"Ауд. 2201 (бывш. 206), Кронверкский пр., д.49,
лит.А",Очно-дистанционный
"сб, 29 октября",13:30,15:00,Лекция,Математика (базовый уровень),Правдин
Константин Владимирович,"Ауд. 1404, Кронверкский пр., д.49, лит.А",Очный
"сб, 29 октября",13:30,15:00,Лекция,Математика (базовый уровень),Правдин
Константин Владимирович,"Ауд. 2201 (бывш. 206), Кронверкский пр., д.49,
лит.А",Очно-дистанционный
"сб, 29 октября",11:40,13:10,Лекция,Математика (базовый уровень),Правдин
Константин Владимирович,"Ауд. 1404, Кронверкский пр., д.49, лит.А",Очный
```

Для решения данного варианта была использована библиотека из задания 2 для создания словаря Python и дальнейшего его парсинга в CSV.

Так для того, чтобы файл в формате XML конвертировать в формат CSV, необходимо, чтобы в файле XML описывалась одинаковая структура объектов, иначе перевод данного формата в CSV окажется бессмысленным.

Вывод

В ходе выполнения лабораторной работы я познакомился с разновидностями форматов обмена информацией и языков разметки, узнал о новых для себя форматах файлов, научился взаимодействовать с ними и реализовал несколько программ, позволяющих переводить файлы в XML формате в форматы YAML и CSV.

Список литературы

- 1) Лямин А.В., Череповская Е.Н. Объектно-ориентированное программирование. Компьютерный практикум. – СПб: Университет ИТМО, 2017. – 143 с. – Режим доступа: <https://books.ifmo.ru/file/pdf/2256.pdf>
- 2) Форма Бэкуса-Наура // Википедия - свободная энциклопедия URL: https://ru.wikipedia.org/wiki/Форма_Бэкуса_—_Наура (дата обращения: 23.10.2022).
- 3) Пишем изящный парсер на Питоне // Хабр URL: <https://habr.com/ru/post/309242/> (дата обращения: 23.10.2022).