

НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ ИТМО

Факультет ФПИИКТ

Дисциплина: Информатика

Лабораторная работа № 4

Выполнил студент

Агнистова Алина Юрьевна

Группа № Р3125

Преподаватель

Болдырева Елена Александровна

г. Санкт-Петербург

2022

Содержание

Задание	3
Отчет	3
Реализация 1 задания:	5
Реализация 2 задания:	5
Реализация 3 задания:	5
Вывод.....	6
Список литературы	7

Задание

1. Обязательное задание

Написать программу на языке Python 3.x, которая бы осуществляла парсинг и конвертацию исходного файла в новый.

2. Доп. задание No1 (+10 баллов)

- а) Найти готовые библиотеки, осуществляющие аналогичный парсинг и конвертацию файлов.
- б) Переписать исходный код, применив найденные библиотеки. Регулярные выражения также нельзя использовать.
- с) Сравнить полученные результаты и объяснить их сходство/различие.

3. Доп. задание No2 (+10 баллов)

- а) Переписать исходный код, добавив в него использование регулярных выражений.
- б) Сравнить полученные результаты и объяснить их сходство/различие.

4. Доп. задание No3 (+10 баллов)

- а) Используя свою исходную программу из обязательного задания, программу из дополнительного задания No1 и программу из дополнительного задания No2, сравнить десятикратное время выполнения парсинга + конвертации в цикле.
- б) Проанализировать полученные результаты и объяснить их сходство/различие.

5. Доп. задание No4 (+5 баллов)

- а) Переписать исходную, чтобы она осуществляла парсинг и конвертацию исходного файла в любой другой формат (кроме JSON, YAML, XML, HTML): PROTOBUF, TSV, CSV, WML и т.п.
- б) Проанализировать полученные результаты, объяснить особенности использованного формата.

Отчет

Реализация основного задания

```
base.py x
1  from time import time #dop3
2  start = time() #dop3 (точка отсчёта времени)
3  def yaml_to_list(filename):
4      with open(filename, "r") as file:
5          strings = file.read().split('\n')
6          dayname = ''
7          count = -1 #counter of the number of lessons
8          main = {} #nice sequence to do json
9          for i in strings:
10             strr = i.lstrip().rstrip()
11             if strr.count(':') == 1:
12                 key, expression = strr.split(':')
13                 if '-' not in key and len(expression) == 0:
14                     main[key] = []
15                     dayname = key
16                 elif '-' in key:
17                     temp = key.lstrip('- ')
18                     main[dayname].append({temp: {}})
19                     count += 1
20                     lessonatm = temp # we save lesson what we have at the moment
21                 else:
22                     expression = expression.lstrip().rstrip()
23                     if '"' in expression:
24                         expression = expression.lstrip('"').rstrip('"')
25                         main[dayname][count][lessonatm][key] = expression
26             elif strr.count(':') == 3:
27                 key, timepart1, timepart2, timepart3 = strr.split(":")
28                 expression = timepart1 + ":" + timepart2 + ":" + timepart3
29                 expression = expression.lstrip().rstrip()
30                 if '"' in expression:
31                     expression = expression.lstrip('"').rstrip('"')
32                 main[dayname][count][lessonatm][key] = expression
33
34  return [main, dayname]
```

```
#####
def list_to_json(main, dayname, output_filename):
    spaces = 2
    with open(output_filename, "w") as file1:
        numb = 0
        file1.write("{\n")
        file1.write(spaces*" " + "'" + dayname + "': [\n")
        spaces += 2
        file1.write(spaces*" " + "{\n")
        spaces += 2
        for i in main[dayname]:
            for key1 in i.keys():
                file1.write(spaces*" " + "'" + key1 + "':{\n")
                temp_main = main[dayname][numb][key1]
                numb += 1
                spaces += 2
                for key2 in temp_main.keys():
                    file1.write(spaces*' ' + "'" + key2 + "': ")
                    if key2 != "teacher":
                        file1.write("'" + temp_main[key2] + "', \n")
                    else:
                        file1.write("'" + temp_main[key2] + '"\n')
                spaces -= 2
                file1.write(spaces*' ' + "]\n")
                if key1 != "ОПД(практика 2)":
                    file1.write(spaces * ' ' + "},\n")
                    file1.write(spaces * ' ' + "{\n")
                else:
                    file1.write(spaces * ' ' + "]\n")
                    spaces -= 2
                    file1.write(spaces * ' ' + "]\n")
                    spaces -= 2
                    file1.write(spaces * ' ' + "]\n")
            spaces += 2
```

```
middle_stage = yaml_to_list("yaml")
last_stage = list_to_json(*middle_stage, "json")
file = open("json").read()
print(file)
print("time base_task:", 10*(time() - start)) #доп3 (время работы программы, время в данный момент - время старта)
```

```

{
  "Пятница": [
    {
      "ОПД(практика 1)": {
        "time": "10:00-11:30",
        "weeks": "2-16",
        "location": "Кронверкский пр., д.49, лит.А",
        "classroom": "2305/1",
        "teacher": "Белозубов Александр Владимирович",
        "type": "Очно-дистанционный",
      }
    },
    {
      "ОПД(практика 2)": {
        "time": "11:40-13:10",
        "weeks": "2-16",
        "location": "Кронверкский пр., д.49, лит.А",
        "classroom": "2305/1",
        "teacher": "Белозубов Александр Владимирович",
        "type": "Очно-дистанционный",
      }
    }
  ]
}

time base_task: 0.006706714630126953

Process finished with exit code 0

```

Реализация 1 доп. задания:

```

dop1.py
1  #with using libs
2  import yaml
3  import json
4  from time import time #dop3
5  sec = time() #dop3 (точка отсчёта времени)
6  def yaml_to_json(input_file, output_file):
7      with open(input_file, 'r', encoding="UTF-8") as yamlfile, open(output_file, "w", encoding="UTF-8") as jsonfile:
8          yaml_dict = yaml.safe_load(yamlfile) #reading file and show it like a dictionary
9          json.dump(yaml_dict, jsonfile, ensure_ascii=False, sort_keys=True, indent=2) #write json file and do a sort dictionaries, indenting, экранирование
10
11  yaml_to_json("yaml", "json1")
12  file = open("json1").read()
13  print(file)
14  print("time dop1_task:", 10*(time() - sec)) #dop3 (время работы программы, время в данный момент - время старта)

```

```

{
  "Пятница": [
    {
      "classroom": "2305/1",
      "location": "Кронверкский пр., д.49, лит.А",
      "teacher": "Белозубов Александр Владимирович",
      "time": "10:00-11:30",
      "type": "Очно-дистанционный",
      "weeks": "2-16",
      "ОПД(практика 1)": null
    },
    {
      "classroom": "2305/1",
      "location": "Кронверкский пр., д.49, лит.А",
      "teacher": "Белозубов Александр Владимирович",
      "time": "11:40-13:10",
      "type": "Очно-дистанционный",
      "weeks": "2-16",
      "ОПД(практика 2)": null
    }
  ]
}

time dop1_task: 0.034432411193847656

Process finished with exit code 0

```

Реализация 2 доп. задания:

```
dop2.py
2 import re
3 from time import time #dop3
4 sec = time() #dop3 (точка отсчёта времени)
5 def from_yaml_to_json_with_regulars(filename, output_filename):
6     with open(filename, 'r', encoding="UTF-8") as yamlf, open(output_filename, "w", encoding="UTF-8") as jsonf:
7         flag_last = False
8         strrr = yamlf.readlines()
9         jsonf.write("{\n")
10        for line in strrr:
11            temp = re.findall(r"^\w+:", line)
12            if len(temp) > 0:
13                result = re.sub(r"^\w+:", " " * 2 + "'" + temp[0].rstrip(":") + "':", line)
14                jsonf.write(result)
15            temp = re.findall(r"^\s+:", line)
16            if len(temp) > 0:
17                if "практика 2" in line:
18                    flag_last = True
19                    jsonf.write("    {\n")
20                    temp = line.lstrip(" ").rstrip(":\n")
21                    jsonf.write(6 * " " + "'" + temp + "': {\n")
22            temp = re.findall(r"^\s+.: .+", line)
23            if len(temp) > 0:
24                temp_2 = temp[0].lstrip().split(":")
25                temp_2[1] = temp_2[1].lstrip(" ").rstrip("'")
26                if len(temp_2) == 3:
27                    jsonf.write(' ' * 8 + "'" + temp_2[0] + "': '" + temp_2[1] + "': '" + temp_2[2].rstrip("'") + "':\n")
28                elif "type" not in temp_2[0]:
29                    jsonf.write(' ' * 8 + "'" + temp_2[0] + "': '" + temp_2[1] + "':\n")
30                else:
31                    jsonf.write(' ' * 8 + "'" + temp_2[0] + "': '" + temp_2[1] + "':\n")
32                    if not flag_last:
33                        jsonf.write(6 * " " + '}\n' + "    },\n")
34                    else:
35                        jsonf.write(6 * " " + '}\n' + "    }\n" + "}")
36        from_yaml_to_json_with_regulars("yaml", "json2")
37        file = open("json2").read()
38        print(file)
39        print("time dop2_task:", 10*(time() - sec)) #dop3 (время работы программы, время в данный момент - время старта)
```

```
{
  "Пятница": [
    {
      "ОПД(практика 1)": {
        "time": "10",
        "weeks": "2-16",
        "location": "Кронверкский пр., д.49, лит.А",
        "classroom": "2305/1",
        "teacher": "Белозубов Александр Владимирович",
        "type": "Очно-дистанционный"
      }
    },
    {
      "ОПД(практика 2)": {
        "time": "11",
        "weeks": "2-16",
        "location": "Кронверкский пр., д.49, лит.А",
        "classroom": "2305/1",
        "teacher": "Белозубов Александр Владимирович",
        "type": "Очно-дистанционный"
      }
    }
  ]
}

time dop2_task: 0.0046539306640625

Process finished with exit code 0
```

Все конвертированные файлы идентичны, за исключением 1 доп. задания и основного, где порядок характеристик другой из-за сортировки. Сравнение времени: основное решение работает быстрее всего, среднее время имеет решение с библиотеками (оно медленнее, чем основное, так как необходимо время на загрузку библиотек), самым медленным оказалось решение с помощью регулярных выражений из-за постоянного вызова библиотеки регулярных выражений.

Вывод

В ходе выполнения я ознакомилась с тем, что представляет из себя парсинг данных и языки разметки, попробовала решить задачу парсинга разными способами, улучшила навык работы с регулярными выражениями.

Список литературы

1. Регулярные выражения в Python от простого к сложному. // Habr URL: <https://habr.com/ru/post/349860/> (дата обращения: 15.11.2022).
2. YAML vs JSON. // Habr URL: https://habr.com/ru/company/rambler_and_co/blog/525498/ (дата обращения: 15.11.2022).