

Работа с разными форматами данных

Николай Свиридов



Николай Свиридов

О спикере:

- Backend-разработчик, IT-блогер



Проверка связи



Отправьте, пожалуйста, смайлик в чат

Если у вас все отлично со связью :-)

Если у вас есть какие-то проблемы со связью :-(



Если у вас нет звука:

- убедитесь, что на вашем устройстве и на колонках включен звук
- обновите страницу вебинара (или закройте страницу и заново присоединитесь к вебинару)
- откройте вебинар в другом браузере
- перезагрузите компьютер (ноутбук) и заново попытайтесь зайти



Поставьте “+”, если меня видно и слышно

План занятия

1. Введение
2. Формат CSV
3. Формат JSON
4. Формат YAML
5. Формат XML
6. Проблема кодировок

ВВЕДЕНИЕ



CSV
плоский

JSON
дерево

YAML
дерево

XML
дерево

Сериализация

Сериализация — процесс преобразования объекта в поток байтов для сохранения или передачи в память, базу данных или файл.

Предназначена для того, чтобы сохранить состояния объекта для последующего воссоздания при необходимости.

Обратный процесс называется **десериализацией**.

FORMAT CSV

Формат CSV

- 1 title,link,description,pubDate
- 2 Как не править Python тесты,https://habr.com/ru/post/502278/?utm_campaign=502278&utm_source=habrahabr&utm_medium=rss, "<p>И вынести тестируемые результаты вне кода. Это статья об автоматизации и увеличения удобства тестирования на Python.</p><p>У меня был проект, который разрабатывался уже несколько лет. В проекте отсутствовали тесты. А также у него были активные зависимости от других команд, которые также влияли на результат.</p><p>Регрессионное тестирование было одним из шагов для более уверенной разработки. Его суть в сравнении вычисленных данных с последним канонизированным результатом работы программы.</p><p>Результаты выполнения можно проверять в python коде тестов. Это близко к контексту выполнения и зачастую удобно.</p>", "Wed, 20 May 2020 11:25:18 +0300"
- 3 Как построить диаграмму Венна с 50 кругами? Визуализация множеств и история моего Python-проекта с открытым кодом,https://habr.com/ru/post/501924/?utm_campaign=501924&utm_source=habrahabr&utm_medium=rss, "<p>Сегодня хочу рассказать вам про задачу

1	title	link	description	pubDate
2	Как не править Python тесты	https://habr.com/ru/post/502278	<p>И вынести тестируемые резуль	20.05.2020 11:25
3	Как построить диаграмму Венна с 50 кругами? Визуализация мно	https://habr.com/ru/post/501924	<p>Сегодня хочу рассказать вам п	20.05.2020 16:13
4	jupyter-book - делаем интерактивную книгу из Jupyter Notebooks	http://github.com/executablebook		19.05.2020 11:14
5	Исключаем дефекты с изображения с помощью OpenCV	https://www.pyimagesearch.com/		20.05.2020 1:38

Формат CSV

Основное применение — **выгрузки данных**

- Для хранения больших объемов (до нескольких Гб) единообразных данных
- Самый компактный формат из всех
- Самый популярный формат для обмена данными у аналитиков
- Поддерживается MS Excel (эквивалентен плоской таблице)
- Не подходит для иерархических данных

<https://docs.python.org/3/library/csv.html>

FORMAT JSON

Формат JSON

```
1  {
2    "channel": {
3      "title": "Дайджест новостей о python",
4      "link": "https://pythondigest.ru/",
5      "description": "Русскоязычные анонсы свежих новостей о python и близлежащих технологиях.",
6      "lastBuildDate": "Wed, 20 May 2020 16:13:18 +0300",
7      "items": [
8        {
9          "title": "Как не править Python тесты",
10         "link": "https://habr.com/ru/post/502278/?utm_campaign=502278&utm_source=habrahabr&utm_medium=rss",
11         "description": "<p>И вынести тестируемые результаты вне кода. Это статья об автоматизации и увеличения удобства тестирования на Python.</p>",
12         "pubDate": "Wed, 20 May 2020 11:25:18 +0300"
13       }
14     ]
15   }
16 }
```

Формат JSON

Основное применение — **базы данных, выгрузки данных**

- Для импорта/экспорта данных в базы данных (в т.ч. bson)
- Для сохранения вложенных структур данных
- При передаче данных клиент <-> сервер для сериализации иерархических объектов
- Самый популярный и простой в использовании формат для Python и Java программистов
- Является подмножеством формата YAML

<https://docs.python.org/3/library/json.html> документация JSON

<https://jsoneditoronline.org/> онлайн редактор JSON

ФОРМАТ YAML

Формат YAML

```
1 channel:
2   description: Русскоязычные анонсы свежих новостей о python и близлежащих технологиях.
3   items:
4     - description: '<p>И вынести тестируемые результаты вне кода. Это статья об автоматизации
5       и увеличения удобства тестирования на Python.</p>'
6       link: https://habr.com/ru/post/502278/?utm_campaign=502278&utm_source=habrahabr&
7         utm_medium=rss
8       pubDate: Wed, 20 May 2020 11:25:18 +0300
9       title: Как не править Python тесты
10 lastBuildDate: Wed, 20 May 2020 16:13:18 +0300
11 link: https://pythondigest.ru/
12 title: Дайджест новостей о python
```

Формат YAML

Основное применение — **файлы конфигурации**

- Самый компактный язык разметки
- Для создания файлов настроек
- Для описания классов, ресурсов и манифестов в API

<https://yaml.org/>

<https://github.com/yaml/pyyaml>

ФОРМАТ XML

XML vs JSON

```
1 <?xml version="1.0" ?>
2 <rss xmlns:ns0="http://www.w3.org/2005/Atom" version="2.0">
3   <channel>
4     <title>Дайджест новостей о python</title>
5     <link>https://pythondigest.ru/</link>
6     <description>Русскоязычные анонсы свежих новостей о python и близлежащих
    технологиях.</description>
7     <ns0:link href="https://pythondigest.ru/rss/" rel="self"/>
8     <language>ru-ru</language>
9     <lastBuildDate>Wed, 20 May 2020 16:13:18 +0300</lastBuildDate>
10    <item>
11      <title>Как не править Python тесты</title>
12      <link>https://habr.com/ru/post/502278/?utm_campaign=502278&utm_source=habrahabr&utm_medium=rss</link>
13      <description>&lt;p&gt;И вынести тестируемые результаты вне кода. Это
        статья об автоматизации и увеличения удобства тестирования на Python.&
        lt;p&gt;</description>
14      <pubDate>Wed, 20 May 2020 11:25:18 +0300</pubDate>
15      <guid>https://habr.com/ru/post/502278/?utm_campaign=502278&utm_source=habrahabr&utm_medium=rss</guid>
16    </item>
17    <item>
18      <title>Как построить диаграмму Венна с 50 кругами? Визуализация
        множеств и история моего Python-проекта с открытым кодом</title>
19      <link>https://habr.com/ru/post/501924/?utm_campaign=501924&utm_source=habrahabr&utm_medium=rss</link>
20      <description>&lt;p&gt;Сегодня хочу рассказать вам про задачу
        визуализации пересекающихся множеств. Поехали!&lt;p&gt;</description>
21      <pubDate>Wed, 20 May 2020 16:13:18 +0300</pubDate>
22      <guid>https://habr.com/ru/post/501924/?utm_campaign=501924&utm_source=habrahabr&utm_medium=rss</guid>
23    </item>
```

```
1 {
2   "channel": {
3     "title": "Дайджест новостей о python",
4     "link": "https://pythondigest.ru/",
5     "description": "Русскоязычные анонсы свежих новостей о ру
    технологиях.",
6     "lastBuildDate": "Wed, 20 May 2020 16:13:18 +0300",
7     "items": [
8       {
9         "title": "Как не править Python тесты",
10        "link": "https://habr.com/ru/post/502278/?utm_campaign=502278&utm_source=habrahabr&utm_medium=rss",
11        "description": "<p&gt;И вынести тестируемые результаты
        статья об автоматизации и увеличения удобства тестир
        Python.</p>",
12        "pubDate": "Wed, 20 May 2020 11:25:18 +0300"
13      },
14      {
15        "title": "Как построить диаграмму Венна с 50 кругами
        множеств и история моего Python-проекта с открытым к
16        "link": "https://habr.com/ru/post/501924/?utm_campaign=501924&utm_source=habrahabr&utm_medium=rss",
17        "description": "<p&gt;Сегодня хочу рассказать вам про з
        пересекающихся множеств. Поехали!</p>",
18        "pubDate": "Wed, 20 May 2020 16:13:18 +0300"
19      }
20    ]
21  }
22 }
```

Формат XML

Основное применение — **сериализация объектов любой сложности**

- Применяется при передаче данных клиент<->сервер для сериализации объектов
- Стандарт обмена данными и сообщениями большинства информационных систем
- Для создания файлов конфигурации

<https://docs.python.org/3/library/xml.etree.elementtree.html>

Формат XML

XPath

Поиск одного элемента: `root.find(query)`

Поиск нескольких элементов: `root.findall(query)`

`query = XPath`

https://www.w3schools.com/xml/xpath_syntax.asp синтаксис XPath

Проблема кодировок

Кодировки: проблемы с кириллицей

Один и тот же текст в кодировке **Windows 1251** и **utf-8**

Windows 1251 (один байт на букву)

00000000:	C8 F1 F2 EE F0 E8 FF 20	F3 20 EC E5 ED FF 20 F1		История у меня с
00000010:	EB E5 E4 F3 FE F9 E0 FF	3A 20 EF EE E7 ED E0 EA		ледующая: позн
00000020:	EE EC E8 EB F1 FF 20 F1	20 E4 E5 E2 F3 F8 EA EE		омился с девушко
00000030:	E9 20 E8 E7 20 F1 E2 EE	E5 E3 EE 20 E8 ED F1 F2		й из своего инст
00000040:	E8 F2 F3 F2 E0 2E 0D 0A	D1 ED E0 F7 E0 EB E0 20		итута...Сначала

utf-8 (два байта на букву)

00000000:	EF BB BF	D0 98 D1 81 D1	82 D0 BE D1 80 D0 B8 D1		п»iP .CfC ,PcCЪPёC
00000010:	8F 20 D1 83 20 D0 BC D0	B5 D0 BD D1 8F 20 D1 81			У Cf PjPµPSCУ Cf
00000020:	D0 BB D0 B5 D0 B4 D1 83	D1 8E D1 89 D0 B0 D1 8F			P»PµPrCfCfC%P°CУ
00000030:	3A 20 D0 BF D0 BE D0 B7	D0 BD D0 B0 D0 BA D0 BE			: PïPpP-PSP°PePs
00000040:	D0 BC D0 B8 D0 BB D1 81	D1 8F 20 D1 81 20 D0 B4			PjPёP»CfCУ Cf Pr

Домашнее задание

Пройдите квиз по теме «Работа с разными форматами данных»



Ваши вопросы?