

# Лабораторная работа 9.

## Создание текстового файла большого объема с помощью генераторной функции

Компьютерная математика II, ММФ, БГУ

Лаврова О.А., идея Козак А.В., Кушнеров А.В., май 2025

### Навыки

1. базовый **Python**: работа с файлами: `with`, `open`, методы файловых объектов `close`, `write`, `writelines`; генераторные выражения и функции; `ord`, `chr`, `len`; методы строковых объектов `join`, `encode`
2. модуль `time` из стандартной библиотеки модулей: `process_time`
3. модуль `random` из стандартной библиотеки модулей: `randint`

### Постановка задачи

Напишите генераторную функцию

```
file_gen(file_name: str,  
         file_size: int,  
         type_symbols: str='latin',  
         number_words: (int,int)=(10,20),  
         number_symbols: (int,int)=(5,10))
```

для создания текстового файла с именем `file_name` размера `file_size` в мегабайтах, строки которого будут состоять из *случайного* количества слов из диапазона `number_words` *случайной* длины `number_symbols` и составленных из *случайных* символов заданного типа `type_symbols`. Слова в строке разделены одним пробелом.

Комментарии по аргументам функции `file_gen`:

- `file_name` определяет имя текстового файла;
- `file_size` определяет размер текстового файла **в мегабайтах**;
- `type_symbols` может принимать одно из трех значений: `'digits'`, `'latin'`, `'cyrillic'`:
  - если `type_symbols='digits'`, то для создания файла используются только цифровые символы;

- если `type_symbols='latin'` , то для создания файла используются только прописные латинские символы от `a` до `z` ; *стандартное значение* `type_symbols='latin'` ;
- если `type_symbols='cyrillic'` , то для создания файла используются только прописные кириллические символы;
- `number_words` является кортежем из двух целых чисел для задания количества слов в строке; количество слов в каждой строке файла определяется как случайное число из диапазона, определенного кортежем; *стандартное значение* `number_words=(10,50)` ;
- `number_symbols` является кортежем из двух целых чисел для задания количества символов в слове; количество символов в каждом слове файла определяется как случайное число из диапазона, определенного кортежем; *стандартное значение* `number_symbols=(5,9)` .

## Функция randint из модуля random

Для выполнения заданий лабораторной работы будем использовать функцию генерации целых чисел `randint` из модуля `random` для создания:

- случайного количества слов в строке файла по заданному диапазону `number_words` ,
- случайной длины слова в строке файл по заданному диапазону `number_symbols` ,
- случайного символа в слове из коллекции `type_symbols` .

```
In [10]: from random import randint
```

```
In [12]: ?randint
```

**Signature:** randint(a, b)

**Docstring:**

Return random integer in range [a, b], including both end points.

**File:** c:\users\user\anaconda3\lib\random.py

**Type:** method

Приведем пример создания символа прописной латинской буквы по его коду, сгенерированному случайно:

```
In [15]: latin_code = (ord('a'), ord('z')) # соответствуем type_symbols = 'Latin'  
f'Код символа: {(c:=randint(*latin_code))}, символ: {chr(c)}'
```

```
Out[15]: 'Код символа: 117, символ: u'
```

## Функция process\_time из модуля time

Для выполнения заданий лабораторной работы будем использовать функцию `process_time` из модуля `time`, которая возвращает значение времени, затраченное процессором на выполнение кода

```
In [19]: from time import process_time
```

```
In [21]: ?process_time
```

**Docstring:**

`process_time()` -> float

Process time for profiling: sum of the kernel and user-space CPU time.

**Type:** builtin\_function\_or\_method

Пример использования функции `process_time`:

```
In [23]: start = process_time()
[x**2 for x in range(1_000_000)]
end = process_time()
print(f'{end - start} секунд')
```

0.203125 секунд

## Задание 9.1. Запись данных в текстовый файл

Откроем файл с именем `test.txt` в текстовом режиме для записи

```
In [27]: file_name = 'test.txt'
f = open(file_name, 'w')
```

Создадим список, состоящий из некоторых строковых объектов, каждый из которых заканчивается символом новой строки `\n`

```
In [30]: lines_list = [f'{x}\n' for x in 'test']
lines_list
```

```
Out[30]: ['t\n', 'e\n', 's\n', 't\n']
```

Запишем созданные строки в файл с помощью метода `writelines` файлового объекта `f`

```
In [33]: f.writelines(lines_list)
```

Вызовем метод `close` файлового объекта `f` для прекращения его связи с внешним файлом, освобождения системных ресурсов и сброса буферизованного вывода на диск, если он находится в памяти.

```
In [36]: f.close()
```

Прочитаем записанные данные. Для этого откроем файл в режиме для чтения `'r'`

```
In [39]: with open(file_name, 'r') as f:
         for line in f:
             print(line, end='')
```

t  
e  
s  
t

В случае записи в файл больших объемов данных целесообразнее использовать не списки для хранения строк, а генераторные объекты

```
In [42]: lines_gen_expr = (f'{x}\n' for x in range(10**6))

def lines_gen_fun(number_lines=10**6):
    yield from (f'{x}\n' for x in range(number_lines))

with open(file_name, 'w') as f:
    f.writelines(lines_gen_expr)
    f.writelines(lines_gen_fun())
```

**В дальнейшем будем создавать данные для записи в файл с помощью генераторной функции.**

Добавим в функцию `lines_gen_fun` подсчет длины записанной информации в мегабайтах. Полагаем, что длина одного символа равна одному байту

```
In [44]: def lines_gen_fun(number_lines = 10**6):
         file_size = 0
         for x in range(number_lines):
             line = f'{x}\n'
             yield line
             file_size += len(line)
         print(f'{file_size/1024**2} Mb')

with open(file_name, 'w') as f:
    f.writelines(lines_gen_fun())
```

6.569757461547852 Mb

**Комментарий к коду.** Допущение, что размер символа в строке равен одному байту, не всегда верно. В частности, для кириллических символов это не так. Более правильным решением для подсчета длины записанной информации в байтах будет использование метода `encode` для записываемых строк. В этом случае в коде выше вместо `len(line)` стоит вызывать `len(line.encode('utf8'))`.

```
In [90]: len('f'), len('f'.encode('utf8')), len('ц'), len('ц'.encode('utf8'))
```

```
Out[90]: (1, 1, 1, 2)
```

Измеряем время, необходимое для создания файла с помощью функции `lines_gen_fun`

```
In [51]: with open(file_name, 'w') as f:
         start = process_time()
         f.writelines(lines_gen_fun())
```

```
end = process_time()
print(end - start, 'секунд')
```

6.569757461547852 Mb  
3.625 секунд

Внесем изменения в генераторную функцию `lines_gen_fun`, чтобы она отображала процент записанных данных в файл, если аргумент `status=True`. Это существенно увеличит время выполнения функции!

```
In [53]: def lines_gen_fun(number_lines=10**6, status=False):
         file_size = 0
         for x in range(number_lines):
             line = f'{x}\n'
             yield line
             file_size += len(line)
             # отображение статуса записи в файл в процентах
             if status:
                 status_number = x/number_lines*100
                 # \r возврат курсора в начало строки для перезаписи
                 print(f'\r{int(status_number)}%', end='', flush=True)
         print(f'\n {file_size/1024**2} Mb')
```

Запишем данные в файл:

```
In [57]: with open(file_name, 'w') as f:
         start = process_time()
         f.writelines(lines_gen_fun(3*10**3, status=True))
         end = process_time()
         print(f'{end - start} секунд')
```

99%  
0.013246536254882812 Mb  
6.328125 секунд

**Напишите** комментарии для каждой строки кода функции `lines_gen_fun`.

## Задание 9.2. Генерация строк

Определим значения переменных, которые задают параметры создания строки для записи в файл

```
In [62]: latin_code = (ord('a'), ord('z')) # соответствуюем type_symbols = 'Latin'
         number_words = (10, 20)
         number_symbols = (5, 10)
```

Создадим одно слово случайной длины из интервала `[5, 10]`, состоящее из случайно выбранных латинских символов

```
In [65]: len_word = randint(*number_symbols)
         word = ''.join([chr(randint(*latin_code)) for i in range(len_word)])
         word
```

Out[65]: 'vgfre'

1. **Напишите код** для генерации строки со случайным количеством слов, случайной длиной слова и случайным набором символов в слове для трех множеств смволов: прописные латинские символы, цифровые символы, прописные кириллические символы. **Протестируйте** написанный код.
2. **Напишите** генераторную функцию `lines_gen_fun_v2(file_size, type_symbols, number_words, number_symbols)` для генерации строк, суммарный размер которых в мегабайтах равен `file_size` .
3. **Напишите** комментарии для каждой строки кода функции `lines_gen_fun_v2` .
4. **Напишите** строки документации для функции `lines_gen_fun_v2` .
5. **Протестируйте** функцию `lines_gen_fun_v2` для различных значений аргументов в предположении, что корректность вводимых данных гарантируется.

## Задание 9.3. Запись сгенерированных строк в текстовый файл

1. **Напишите** результирующую генераторную функцию

```
file_gen(file_name: str,  
         file_size: int,  
         type_symbols: str='latin',  
         number_words: (int,int)=(10,20),  
         number_symbols: (int,int)=(5,10))
```

Функция `file_gen` должна использовать функцию `lines_gen_fun_v2` из Задания 9.2.

Функция `file_gen` должна выводить на экран:

- процент записанных данных в процессе выполнения функции,
- фактический размер записанных данных в мегабайтах,
- время для выполнения кода функции.

2. **Напишите** комментарии для каждой строки кода функции `file_gen` .
3. **Напишите** строки документации для функции `file_gen` .
4. **Протестируйте** функцию `file_gen` для различных значений аргументов в предположении, что корректность вводимых данных гарантируется.