

Цель работы.

Обработка набора данных средствами Python.

Задание.

1. С помощью примера кода сгенерировать CSV-файл с данными в 5 столбцов, количество записей 10^6 .
2. Средствами Python прочитать содержимое CSV-файла.
3. Найти среднее значение по колонке 2.
4. Построить график значений для колонок 3 и 4 в зависимости от данных времени из колонки 1.
5. Для колонки 5 построить график медиан кортежей по 10^5 записей. По оси X откладываем номер кортежа.

Параметры для генерации функций приведены в таблице 1.

Таблица 1

Вариант	period	v1	v2	v3
1	100000	3000	1000	200.0
2	130000	2000	500	170.0
3	170000	3000	800	150.0
4	200000	4000	1100	130.0
5	110000	5000	1200	220.0
6	220000	2500	1000	210.0
7	400000	3000	1300	190.0

Теоретические положения.

В качестве основы использован пример кода из задания на языке программирования Python.

При выполнении задания используется типовой подход к моделированию и обработке временных рядов. Исходные данные формируются в виде CSV-файла, содержащего столбец времени и четыре вычисляемые функции. Параметры генерации (period, v1, v2, v3) выбираются согласно таблице 1 для соответствующего варианта.

Функции представляют собой комбинацию детерминированной составляющей и случайного шума. В частности, для func1 используется равномерная случайная добавка к базовому уровню; func2 задаётся синусоидальной зависимостью от времени с шумовой компонентой; func3 формируется как сумма func1 и func2; func4 дополнительно включает аргумент времени, поэтому имеет выраженный тренд роста

по оси времени. Такая структура данных позволяет проверить чтение/запись CSV, вычисление статистик и построение графиков.

При обработке данных выполняются: чтение файла построчно, расчёт среднего значения по второй колонке (func1), построение графиков func2 и func3 от времени, а также расчёт медиан func4 по блокам объёмом 10^5 записей. Медиана используется как устойчивая характеристика центральной тенденции при наличии случайной составляющей.

Ход работы:

Работа выполнена по варианту 6. Для удобства выполнения и проверки задания программа разбита на несколько отдельных скриптов: генерация исходных данных в CSV, вычисление среднего значения, построение графика для колонок 3 и 4 и построение графика медиан для колонки 5. Такое разбиение позволяет запускать этапы независимо и проще контролировать результат каждого шага.

1. Формирование исходного набора данных (CSV)

На первом этапе был сформирован файл test.csv, содержащий 5 столбцов и 10^6 строк. Генерация выполнялась скриптом Генерация CSV.py. В качестве времени используется целочисленная последовательность, начинающаяся с текущей метки времени UNIX (timestamp), что обеспечивает монотонный рост аргумента времени и удобную привязку данных к оси X при построении графиков. Инициализация времени и запуск генерации выполняются следующим образом:

```
1. start = int(datetime.datetime.now().timestamp())
2. generate_csv('test.csv', start, 10 ** 6, period, v1, v2, v3)
```

Параметры функций взяты из таблицы вариантов (вариант 6) и заданы явно в основной части скрипта:

```
1. period = 220000
2. v1 = 2500
3. v2 = 1000
4. v3 = 210.0
```

Структура CSV фиксирована: time, func1, func2, func3, func4. Запись выполняется через csv.DictWriter, что позволяет формировать строки по именам столбцов и

избегать ошибок с порядком колонок. В цикле по i формируются значения функций и добавляются в файл:

```
1. fieldnames = ['time', 'func1', 'func2', 'func3', 'func4']
2. ...
3. writer.writerow({
4.     'time': start + i,
5.     'func1': func_rand(start + i),
6.     'func2': func_sin(start + i, period, v1, v2, v3),
7.     'func3': func_three(start + i, period, v1, v2, v3),
8.     'func4': func_four(start + i, period, v1, v2, v3)
9. })
```

Содержательно столбцы формируются так:

- `func1` — случайная составляющая (базовый уровень + равномерная добавка), что задаёт разброс значений:

```
1. res = 2500 + random.random() * 1000
```

- `func2` — синусоидальная зависимость от времени с шумом. В формуле присутствуют параметры варианта (`period`, `v1`, `v2`, `v3`), где `v1` задаёт уровень, `v2` — амплитуду синусоидальной части, `v3` — уровень случайной добавки:

```
1. s = math.sin((2 * np.pi / period) * time)
2. res = v1 + s * v2 + random.random() * v3
```

По итогу этапа сформирован файл `test.csv`, который далее используется во всех последующих пунктах задания. Блок схема алгоритма генерации представлена на Рисунке 1.

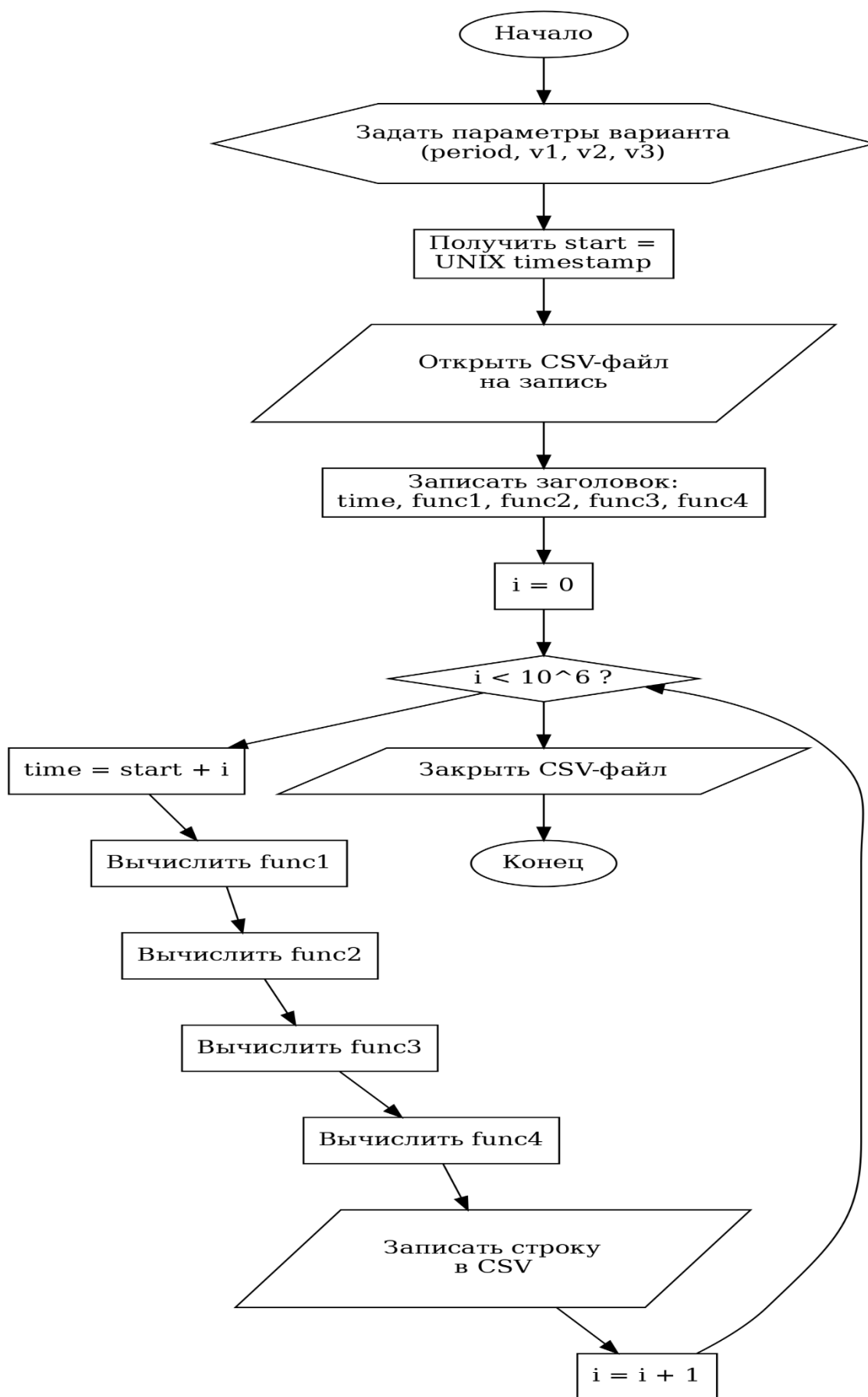


Рисунок 1. Блок-схема алгоритма генерации CSV-файла (формирование 10^6 записей со столбцами time, func1–func4).

2. Чтение CSV и вычисление среднего значения по колонке 2

Далее выполнено чтение сформированного файла средствами Python и вычислено среднее значение по **второй колонке**. В созданном CSV после time второй колонкой является func1, поэтому среднее считается по полю func1. Для этого используется скрипт Среднее по колонке 2.py.

Чтение организовано последовательно (построчно) через csv.DictReader, что удобно при большом объёме данных и не требует загрузки всего файла в память.

В цикле суммируются значения и считается количество строк:

```
1. for row in reader:
2.     total += float(row['func1'])
3.     count += 1
```

После прохода по файлу среднее вычисляется как отношение суммы к числу элементов:

```
1. return total / count
```

Результат выводится в консоль, что позволяет быстро проверить корректность чтения файла и получить численное значение среднего:

```
1. print("Среднее по колонке func1:", avg)
```

3. Построение графика значений для колонок 3 и 4 по времени

Следующий пункт задания — построить графики значений колонок 3 и 4 в зависимости от данных времени (колонка 1). В CSV колонке 3 соответствует func2, колонке 4 — func3. Построение выполнено скриптом График значений для колонок 3 и 4.py.

Сначала выполняется чтение файла и формирование массивов значений по времени и двум функциям:

```
1. times.append(float(row['time']))
2. col3.append(float(row['func2']))
3. col4.append(float(row['func3']))
```

Затем строятся две кривые на одном графике, для визуального сравнения поведения func2 и func3 на одном интервале времени:

```
1. plt.plot(times, col3, label='func2 (колонка 3)')
2. plt.plot(times, col4, label='func3 (колонка 4)')
```

Для корректного оформления добавлены подписи осей, заголовок, легенда и сетка. График сохраняется в PNG с разрешением 300 dpi:

```
1. plt.savefig("График значений для колонок 3 и 4.png", dpi=300)
```

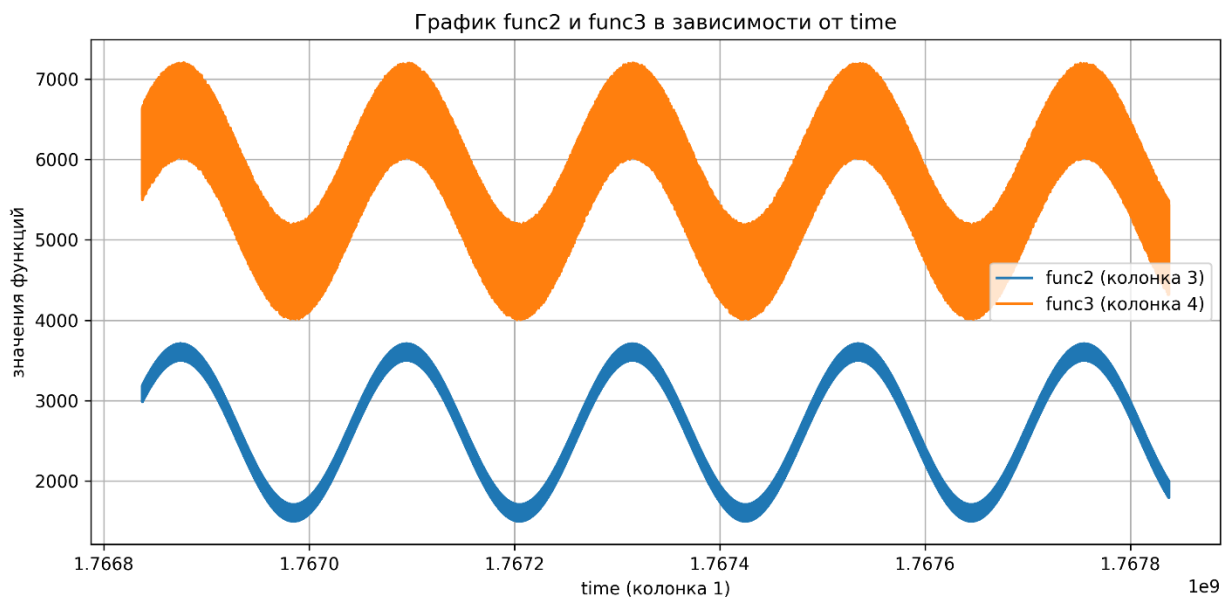


Рисунок 2 — График функций func2 и func3 в зависимости от времени.

4. Построение графика медиан для колонки 5 по блокам 10^5

Заключительный пункт задания — для колонки 5 (то есть func4) построить график медиан по кортежам (блокам) объёмом 10^5 записей, откладывая по оси X номер кортежа. Для этого используется скрипт Строим график медиан.ру. Размер блока задан как 100 000, а число блоков — 10, что соответствует общему объёму 10^6 записей:

```
1. block_size = 100_000
2. blocks = [[] for _ in range(10)]
```

При чтении файла значения func4 распределяются по блокам по индексу строки:

```
1. value = float(row['func4'])
2. block_index = index // block_size
3. if block_index < 10:
4.     blocks[block_index].append(value)
5. index += 1
```

После заполнения блоков для каждого рассчитывается медиана стандартной функцией statistics.median, затем строится график медиан по номерам блоков 1...10:

```
1. medians = [statistics.median(block) for block in blocks]
2. ...
3. plt.plot(range(1, 11), medians, marker='o')
4.
```

График оформляется подписями и сохраняется в PNG:

plt.savefig ("График медиан func4.png", dpi=300)

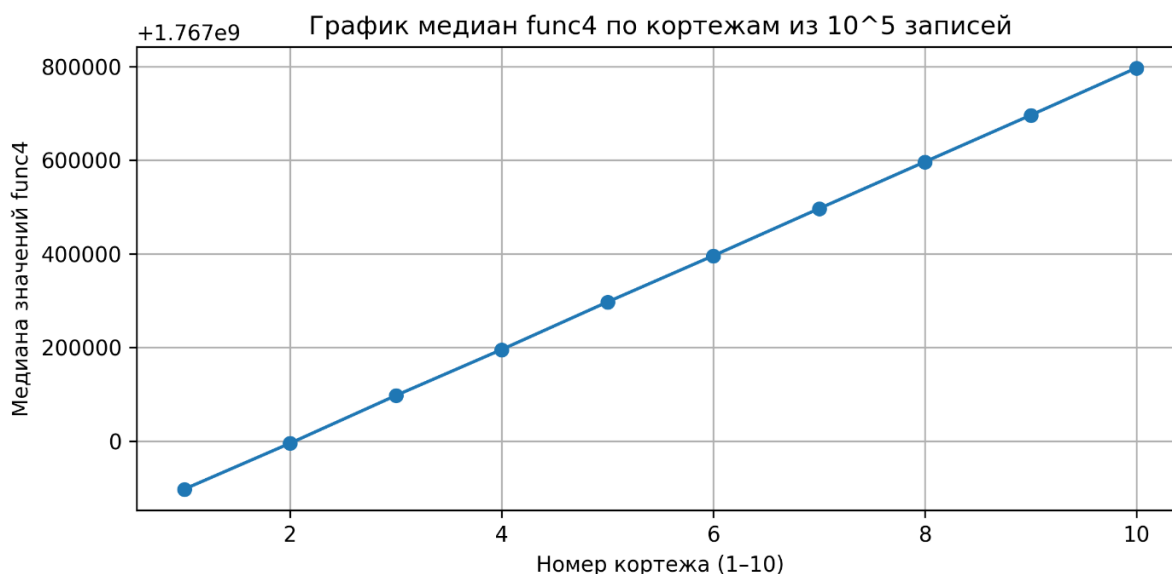


Рисунок 3 — График медиан func4 по блокам из 10^5 записей.

В результате выполнены все пункты задания: сгенерирован CSV-файл заданного формата и объёма, выполнено чтение данных, рассчитано среднее по второй колонке (func1), построен график зависимостей для колонок 3 и 4 от времени, а также построен график медиан для колонки 5 по блокам 10^5 записей. Полученные изображения графиков сохранены в файлы PNG.

Вывод:

В рамках выполнения ИДЗ №1 (вариант 6) были реализованы и последовательно выполнены все пункты задания по обработке набора данных средствами Python. На основе заданных параметров сформирован CSV-файл test.csv, содержащий 10^6 записей и 5 столбцов (временная шкала и четыре функции). После генерации выполнено чтение файла и получено среднее значение по второй колонке (func1), что подтверждает корректность структуры CSV и пригодность данных для дальнейшей обработки.

Для анализа поведения функций построены графики значений func2 и func3 в зависимости от времени, что позволяет наглядно сравнить характер изменения данных в третьей и четвёртой колонках. Также для пятой колонки (func4) выполнено разбиение на 10 блоков по 10^5 значений и построен график медиан по каждому блоку. Такой подход даёт обобщённую оценку центрального уровня значений на больших объёмах данных и позволяет сравнить блоки между собой. Все графики сохранены в формате PNG и использованы в отчёте как иллюстрации результатов. Код модулей представлен в приложениях к отчету.

Таким образом, поставленная цель работы — обработка набора данных средствами Python — достигнута, а требуемые вычисления и визуализация выполнены в полном объёме.

Приложения.

Код исполняемого файла Генерация CSV.py

```
1. import datetime
2. import csv
3. import math
4. import random
5. import numpy as np
6.
7.
8. def func_rand(time:float):
9.     # ""Генерирует случайное значение для столбца func1""
10.    res = 2500 + random.random() * 1000
11.    return res
12.
13. def func_sin(time:float, period:int = 220000, v1:int = 2500, v2:float = 1000, v3:float = 210.0):
14.    # ""Синусоидальная функция с шумом для столбца func2""
15.    s = math.sin((2 * np.pi / period) * time)
16.    res = v1 + s * v2 + random.random() * v3
17.    return res
18.
19. def func_three(time:float, period:int = 220000, v1:int = 2500, v2:float = 1000, v3:float = 210.0):
20.    # ""Комбинация func_sin и func_rand для столбца func3""
21.    res = func_sin(time, period, v1, v2, v3) + func_rand(time)
22.    return res
23.
24. def func_four(time:float, period:int = 220000, v1:int = 2500, v2:float = 1000, v3:float = 210.0):
25.    # ""time + func_three + случайный шум — данные для столбца func4""
26.    res = time + func_three(time, period, v1, v2, v3) + func_rand(time)
27.    return res
28.
29.
30. def generate_csv(filename:str, start:int, count:int, period:int, v1:int, v2:float, v3:float):
31.    # ""Генерирует CSV-файл с пятью функциями по заданным параметрам""
32.    with open(filename, 'w', newline='') as csvfile:
33.        fieldnames = ['time', 'func1', 'func2', 'func3', 'func4']
34.        writer = csv.DictWriter(csvfile, fieldnames=fieldnames)
35.        writer.writeheader()
36.        for i in range(count):
37.            writer.writerow({
38.                'time': start + i,
39.                'func1': func_rand(start + i),
40.                'func2': func_sin(start + i, period, v1, v2, v3),
41.                'func3': func_three(start + i, period, v1, v2, v3),
42.                'func4': func_four(start + i, period, v1, v2, v3)
43.            })
44.
45.
46. if __name__ == '__main__':
47.    start = int(datetime.datetime.now().timestamp())
48.
49.    # параметры варианта 6
50.    period = 220000
51.    v1 = 2500
52.    v2 = 1000
53.    v3 = 210.0
54.
55.    generate_csv('test.csv', start, 10 ** 6, period, v1, v2, v3)
56.    print("CSV-файл 'test.csv' успешно сгенерирован.")
```

Код исполняемого файла График значений для колонок 3 и 4.py

```
1. import csv
2. import matplotlib.pyplot as plt
3.
4. def plot_columns(filename: str):
5.     times = []
6.     col3 = [] # func2
7.     col4 = [] # func3
8.
9.     with open(filename, 'r') as f:
10.         reader = csv.DictReader(f)
11.         for row in reader:
12.             times.append(float(row['time']))
13.             col3.append(float(row['func2']))
14.             col4.append(float(row['func3']))
15.
16.     # Строим графики
17.     plt.figure(figsize=(10, 5))
18.     plt.plot(times, col3, label='func2 (колонка 3)')
19.     plt.plot(times, col4, label='func3 (колонка 4)')
20.
21.     plt.xlabel('time (колонка 1)')
22.     plt.ylabel('значения функций')
23.     plt.title('График func2 и func3 в зависимости от time')
24.     plt.legend()
25.     plt.grid(True)
26.     plt.tight_layout()
27.     plt.tight_layout()
28.     plt.savefig("График значений для колонок 3 и 4.png", dpi=300)
29.     plt.show()
30.
31.
32. if __name__ == '__main__':
33.     plot_columns('test.csv')
34.     print("График успешно сохранён ")
```

Код исполняемого файла Среднее по колонке 2.py

```
1. import csv
2.
3. def read_and_avg(filename: str):
4.     total = 0.0
5.     count = 0
6.
7.     with open(filename, 'r') as f:
8.         reader = csv.DictReader(f)
9.
10.        for row in reader:
11.            total += float(row['func1'])
12.            count += 1
13.
14.        return total / count
15.
16.
17. if __name__ == '__main__':
18.     avg = read_and_avg('test.csv')
19.     print("Среднее по колонке func1:", avg)
20.
```

Код исполняемого файла Строим график медиан.py

```
1. import csv
2. import statistics
3. import matplotlib.pyplot as plt
4.
5. def plot_medians(filename: str):
6.     block_size = 100_000      # 10^5 записей
7.     blocks = [] for _ in range(10) # 10 кортежей
8.
9.     with open(filename, 'r') as f:
10.         reader = csv.DictReader(f)
11.
12.         index = 0
13.         for row in reader:
14.             value = float(row['func4'])
15.             block_index = index // block_size
16.             if block_index < 10:
17.                 blocks[block_index].append(value)
18.             index += 1
19.
20.     # вычисляем медианы каждого блока
21.     medians = [statistics.median(block) for block in blocks]
22.
23.     # строим график
24.     plt.figure(figsize=(8, 4))
25.     plt.plot(range(1, 11), medians, marker='o')
26.     plt.xlabel("Номер кортежа (1–10)")
27.     plt.ylabel("Медиана значений func4")
28.     plt.title("График медиан func4 по кортежам из 10^5 записей")
29.     plt.grid(True)
30.     plt.tight_layout()
31.     plt.tight_layout()
32.     plt.savefig("График медиан func4.png", dpi=300)
33.     plt.show()
34.
35.
36. if __name__ == '__main__':
37.     plot_medians('test.csv')
38.     print("График медиан func4 успешно построен.")
```