

Only you can see this message



This story's distribution setting is on. [Learn more](#)

Python 3.9 StatsProfile — My first OSS Contribution to cPython



Daniel Olshansky

Feb 17 · 4 min read

You can try out all of the code in this article yourself using this [Google Colaboratory notebook](#).

If you've ever tried to debug and optimize your python application, it's likely that you stumbled upon Python Profiles to understand where most of the execution time is being spent. You enable the profiler at the beginning of a code segment you're interested in profiling with `pr.enable()`, and call `pr.create_stats()` at the end. *

```
1  import cProfile, pstats
2  import time
3  from random import randint
4
5  def sleep1():
6      time.sleep(0.1)
7
8  def sleep2():
9      time.sleep(0.2)
10
11  def sleep3():
12      time.sleep(0.3)
13
14  pr = cProfile.Profile()
15  pr.enable()
```

```

15 pr.enable()
16
17 for _ in range(10):
18     for _ in range(randint(1, 10)):
19         sleep1()
20
21     for _ in range(randint(1, 10)):
22         sleep2()
23
24     for _ in range(randint(1, 10)):
25         sleep3()
26
27 pr.create_stats()
28 ps = pstats.Stats(pr)
29 ps.print_stats()

```

get state profile example, print state by hosted with ❤️ by GitHub

[view raw](#)

Afterwards, you can create a Stats object, and print the results in a human readable format with `ps.print_stats()`.

```

66 function calls in 4.507 seconds

Ordered by: internal time, cumulative time

ncalls  tottime  percall  cumtime  percall  filename:lineno(function)
    24     4.507     0.188     4.507     0.188 {built-in method time.sleep}
    10     0.000     0.000     1.002     0.100 <ipython-input-7-96f3d7189345>:12(sleep1)
     7     0.000     0.000     1.402     0.200 <ipython-input-7-96f3d7189345>:15(sleep2)
     7     0.000     0.000     2.103     0.300 <ipython-input-7-96f3d7189345>:18(sleep3)
     3     0.000     0.000     0.000     0.000 /usr/lib/python3.6/random.py:173(randrange)
     3     0.000     0.000     0.000     0.000 /usr/lib/python3.6/random.py:223(_randbelow)
     3     0.000     0.000     0.000     0.000 /usr/lib/python3.6/random.py:217(randint)
     1     0.000     0.000     0.000     0.000 /usr/lib/python3.6/cProfile.py:50(create_stats)
     4     0.000     0.000     0.000     0.000 {method 'getrandbits' of '_random.Random' objects}
     3     0.000     0.000     0.000     0.000 {method 'bit_length' of 'int' objects}
     1     0.000     0.000     0.000     0.000 {method 'disable' of '_lsprof.Profiler' objects}

<pstats.Stats at 0x7f5d60f80a20>

```

The output above is quite useful and can take you a long way. However, what if you don't know what kind of data inputs cause a bottleneck in your application? What if you're interested in aggregating and evaluating profiling data over some period of time? What if you want to profile your application while your team is dogfooding it? I found that there isn't an easy way to use this data in an ETL pipeline where you'd be able to do further offline analysis over a larger dataset.

• • •

I recently made my first open source cPython contribution which adds a dataclass called `StatsProfile` to address this. After you created your stats object, you can retrieve all the information in the above screenshot by calling `ps.get_stats_profile()` and analyze it in a programmatic way.

If you're not on Python3.9 yet, the following code snippet is a slightly modified version of the code in the pull request that you can start using today by importing it directly into your project.

```

1  from pstats import func_std_string, f8
2  from dataclasses import dataclass
3  from typing import Dict
4
5  @dataclass(unsafe_hash=True)
6  class FunctionProfile:
7      ncalls: int
8      tottime: float
9      percall_tottime: float
10     cumtime: float
11     percall_cumtime: float
12     file_name: str
13     line_number: int
14
15  @dataclass(unsafe_hash=True)
16  class StatsProfile:
17     '''Class for keeping track of an item in inventory.'''
18     total_tt: float
19     func_profiles: Dict[str, FunctionProfile]
20
21  def get_stats_profile(stats):
22     """This method returns an instance of StatsProfile, which contains a mapping
23     of function names to instances of FunctionProfile. Each FunctionProfile
24     instance holds information related to the function's profile such as how
25     long the function took to run, how many times it was called, etc...
26     """
27     func_list = stats.fcn_list[:] if stats.fcn_list else list(stats.stats.keys())
28     if not func_list:
29         return StatsProfile(0, {})
30

```

```

30
31     total_tt = float(f8(stats.total_tt))
32     func_profiles = {}
33     stats_profile = StatsProfile(total_tt, func_profiles)
34
35     for func in func_list:
36         cc, nc, tt, ct, callers = stats.stats[func]
37         file_name, line_number, func_name = func
38         ncalls = str(nc) if nc == cc else (str(nc) + '/' + str(cc))
39         tottime = float(f8(tt))
40         percall_tottime = -1 if nc == 0 else float(f8(tt/nc))
41         cumtime = float(f8(ct))
42         percall_cumtime = -1 if cc == 0 else float(f8(ct/cc))
43         func_profile = FunctionProfile(
44             ncalls,
45             tottime, # time spent in this function alone
46             percall_tottime,
47             cumtime, # time spent in the function plus all functions that this function
48             percall_cumtime,
49             file_name,
50             line_number
51         )
52         func_profiles[func_name] = func_profile
53
54     return stats_profile

```

Now, rather than inspecting the profile of a single execution of our code snippet, we can aggregate and analyze the profiles over several different iterations. In a real production service, depending on which logging tool you use, you would likely need to format and stringify the `StatsProfile` dataclass before logging it, but for the purposes of this example, everything is stored in memory.

To simulate timestamped logging, `(timestamp, stats_profile)` tuples are appended to a `timestamped_stats_profile` list with every execution of the loop.

```

1  import cProfile, pstats
2  import time
3  from random import randint
4
5  START_TIME = int(time.time())
6
7  timestamped_stats_profiles = []

```

```
8
9  def sleep1():
10     time.sleep(0.1)
11
12  def sleep2():
13     time.sleep(0.2)
14
15  def sleep3():
16     time.sleep(0.3)
17
18  for _ in range(10):
19     pr = cProfile.Profile()
20     pr.enable()
21
22     for _ in range(randint(1, 10)):
23         sleep1()
24
25     for _ in range(randint(1, 10)):
26         sleep2()
27
28     for _ in range(randint(1, 10)):
29         sleep3()
30
31     pr.create_stats()
32     ps = pstats.Stats(pr)
33
34     stats_profile = get_stats_profile(ps)
35     timestamped_stats_profiles.append((int(time.time()), stats_profile))
```

After the data is logged, it needs to be aggregated over a certain timeslice. Most logging/visualization platforms have their functions to process timeseries data, so this would be platform specific. Sumologic has the timeslice function, Elasticsearch has examples of how to do date histogram aggregation, Datadog has an aggregate across time dropdown, etc...

For the purposes of this example, I'm doing the aggregation manually in python. I bucket all the logged (i.e. saved) StatsProfile objects over 10 second intervals, aggregate the cumulative execution time, `cumtime`, per function call and store the resultant counters in `time_slices_counters`. If you're interested in inspecting the number of calls to certain functions rather than the cumulative execution time spent in

it, you would simply modify the parameter being access on line 21 in the code snippet below.

```

1  from collections import Counter
2  import itertools
3
4  TIME_SLICE = 10 # Aggregate logs every 10 seconds
5
6  def time_to_bucket(time):
7      return (time-START_TIME) // TIME_SLICE
8
9  def bucket_to_time(bucket):
10     return bucket * TIME_SLICE + START_TIME
11
12  time_sliced_counters = []
13  headers = set()
14
15  for bucket, grp in itertools.groupby(timestamped_stats_profiles, key=lambda timestamped_s
16     time_slice_counter = Counter()
17     for (timestamp, stats_profile) in grp:
18         for f_name, f_profile in stats_profile.func_profiles.items():
19             f_name = f_name.lstrip("_") # matplotlib can't allow legend values to start
20             headers.add(f_name)
21             time_slice_counter[f_name] += f_profile.cumtime
22     time_sliced_counters.append((bucket_to_time(bucket), time_slice_counter))

```

get stats profile log aggregation.py hosted with ❤ by GitHub

[view raw](#)

In my opinion, a stacked bar graph is a great way to visualize and easily interpret this data. Using the following code snippet:

```

1  import numpy as np
2  import matplotlib
3  import matplotlib.pyplot as plt
4  from matplotlib.ticker import FormatStrFormatter
5
6  WIDTH = 0.4
7
8  ind = np.arange(len(time_sliced_counters))
9
10 x_axis = tuple(time for (time, c) in time_sliced_counters)
11 y_axis = [[] for _ in range(len(headers))]
12 for idx, header in enumerate(headers):

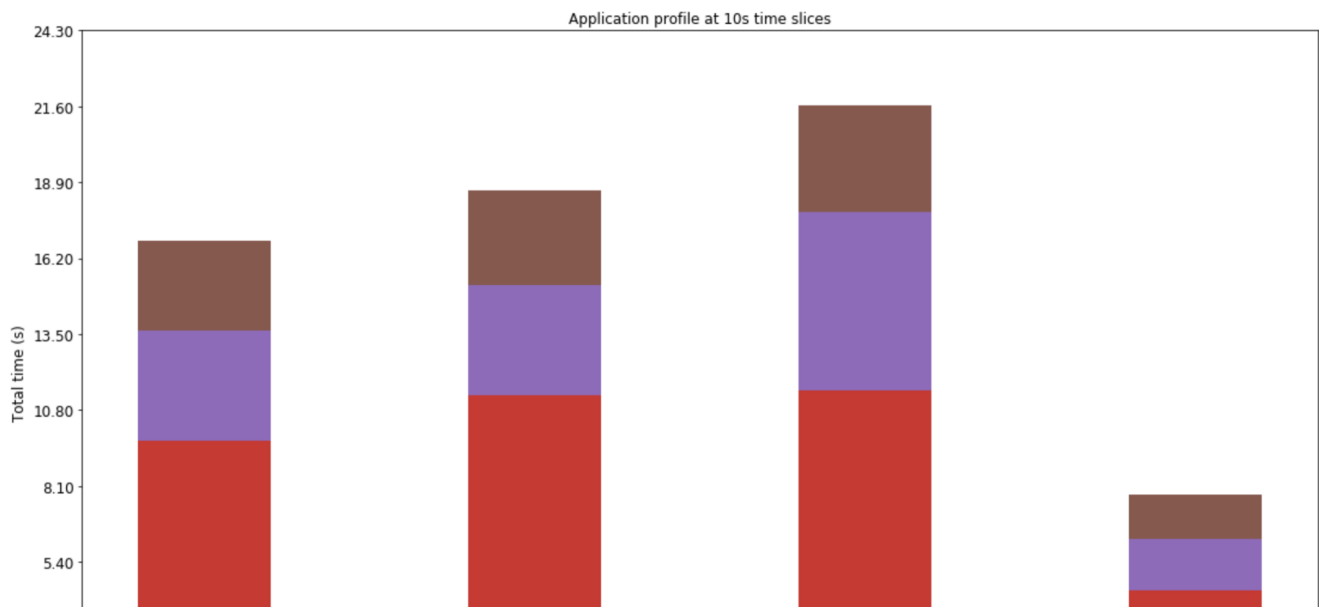
```

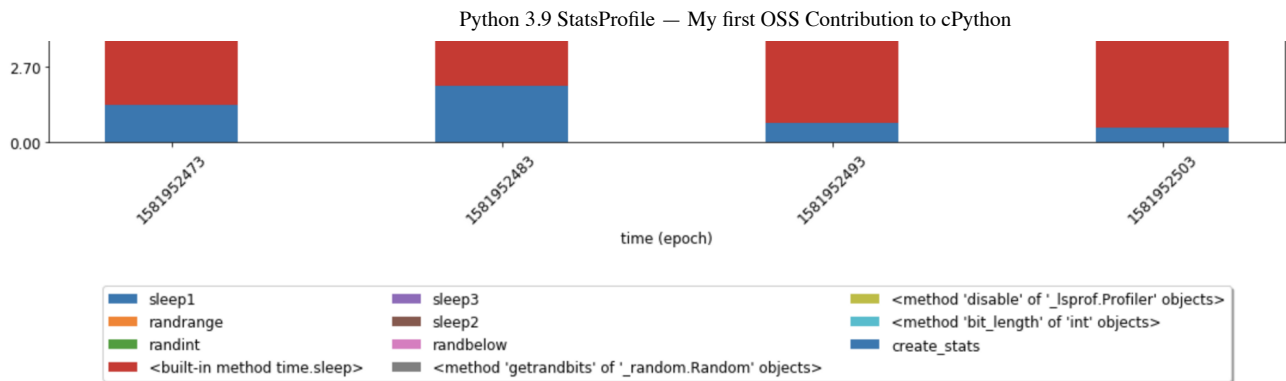
```

13     for (time, counter) in time_sliced_counters:
14         y_axis[idx].append(counter[header])
15
16 fig = matplotlib.pyplot.gcf()
17 fig.set_size_inches(18.5, 10.5)
18
19 boxes = []
20 titles = []
21 bottom = np.zeros(len(y_axis[0]))
22 for idx, header in enumerate(headers):
23     p = plt.bar(ind, tuple(y_axis[idx]), WIDTH, bottom=bottom)
24     bottom += np.array(y_axis[idx])
25     boxes.append(p[0])
26     titles.append(header)
27
28 plt.ylabel("Total time (s)", fontsize=12)
29 plt.yticks(np.arange(0, round(max(bottom) + 5), 1), fontsize=12)
30 plt.gca().yaxis.set_major_formatter(FormatStrFormatter('%.2f'))
31
32 plt.xlabel('time (epoch)', fontsize=12)
33 plt.xticks(ind, x_axis, fontsize=12, rotation=45)
34
35 plt.title(f"Application profile at {TIME_SLICE}s time slices")
36 plt.legend(tuple(boxes), tuple(titles), fontsize=12, ncol=3, fancybox=True, loc='upper c
37
38 plt.show()

```

We can generate a graph that looks like this:





The results aren't very interesting or surprising given the simplicity of a script calling `sleep` a bunch of times, but hopefully it'll be more useful in more complex applications.

• • •

It's important to note that you probably should not be doing this in production. It could be useful on your local or development environments, and might be worth enabling in a single canary, but could have adverse effects in prod. You would be polluting your logs with large `StatsProfile` structures, and I have not investigated if running `cProfile` in prod could potentially downgrade your service's performance.

• • •

As a side note, though there is some overhead and a small learning curve, I was very pleased with how easy it is to contribute to cPython. Aside from publishing the actual PR, you have to sign the PSF Contributor Agreement, open a on bugs.python.org, and nudge a few people to make get your code looked at. There is a great developer guide on how to run things locally and execute tests. I recently also came by this doc, which is a good starting point if you've never contributed to cPython before.

Huge thanks [Gregory P. Smith](#) for reviewing and approving my cPython PR! Also, thank you to [Сергей Яркин](#) for proofreading my article, and a special shoutout to [Manuel Dell'Elce](#) who built a really nice chrome extension that made embedding code snippets in this medium article a breeze.

[Python](#) [Profiling](#) [Stacked Bar Chart](#) [Stats](#)

[About](#) [Help](#) [Legal](#)

Get the Medium app

