

"Online" Algorithms (Data Streams): Ideas and code

Pietro Colaguri 1936709

November 2023

1 Ideas

Online algorithms are algorithms that process data in a sequential order, making decisions without having the entire dataset available at once. In the context of data streams, where data is continuously generated and processed in real-time, online algorithms play a crucial role. These algorithms are designed to handle large datasets that cannot fit into memory and must be processed on the fly.

One common application of online algorithms is in the realm of data stream processing, where the algorithm must process data as it arrives and make decisions based on the current information. Online algorithms are memory-efficient and have a constant or logarithmic space complexity, making them suitable for handling massive streams of data.

I will proceed by discussing the concept of online algorithms in the context of data streams, later I will also provide snippets of code.

2 Code

1. **Number of Elements:** we keep track of all the different elements seen by using a hash table or a probabilistic data structure, like Count-Min Sketch. Below I will show an example script in which I implement this idea using the DS Count-Min Sketch, the output, as expected, will be "Estimated count of item 1: 3".

```
1      from collections import defaultdict
2
3      class CountMinSketch:
4          """
5              Count-Min Sketch is a probabilistic data
6              structure used for estimating the
              frequency of elements in a data stream.
              It uses a fixed amount of memory and
              provides approximate counts with a
              small error rate.
```

```

7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

    Attributes:
        width (int): The width of the sketch
            table.
        depth (int): The depth of the sketch
            table.
        sketch (list): The sketch table used
            for counting.

    Methods:
        hash_func(x, i): Hashes the element x
            using the i-th hash function.
        update(x): Updates the sketch table
            by incrementing the count for
            element x.
        estimate(x): Estimates the frequency
            of element x by finding the
            minimum count across all hash
            functions.

    Usage:
        sketch = CountMinSketch(width, depth)
        sketch.update(x)
        estimated_count = sketch.estimate(x)
    """

def __init__(self, width, depth):
    """
    Initializes a CountMinSketch object.

    Args:
        width (int): The width of the
            sketch table.
        depth (int): The depth of the
            sketch table.
    """
    self.width = width
    self.depth = depth
    self.sketch = [[0] * width for _ in
        range(depth)]

def hash_func(self, x, i):
    """
    Hashes the element x using the i-th
    hash function.

```

```

40         Args:
41             x: The element to be hashed.
42             i (int): The index of the hash
43                 function.
44
45         Returns:
46             int: The hashed value of x.
47         """
48         return hash((x, i)) % self.width
49
50     def update(self, x):
51         """
52         Updates the sketch table by
53         incrementing the count for element
54         x.
55
56         Args:
57             x: The element to be updated.
58         """
59         for i in range(self.depth):
60             j = self.hash_func(x, i)
61             self.sketch[i][j] += 1
62
63     def estimate(self, x):
64         """
65         Estimates the frequency of element x
66         by finding the minimum count
67         across all hash functions.
68
69         Args:
70             x: The element to be estimated.
71
72         Returns:
73             int: The estimated frequency of
74                 element x.
75         """
76         min_count = float('inf')
77         for i in range(self.depth):
78             j = self.hash_func(x, i)
79             min_count = min(min_count,
80                             self.sketch[i][j])
81         return min_count
82
83     # Example usage
84     stream = [1, 2, 3, 1, 2, 4, 5, 1, 2, 3]
85     cms = CountMinSketch(width=100, depth=5)

```

```

79     for item in stream:
80         cms.update(item)
81         count = cms.estimate(1)
82         print(f"Estimated count of item 1: {count}")

```

In order to make the script as clear as possible I added varbose comments.

2. **Sliding Window:** Keep track of a fixed size window containing the latest elements of the data stream, this idea is very powerful and has many application, for instance in Telecommunications with the implementation of a packet-switched network, in which it is needed to keep track inside the buffers of the routers of a fixed size window of packets. To create an interesting script it would be better to implement some sort of data processing on the contents of the window, in this case I decided to compute the average value of the sliding window.

```

1     class SlidingWindow:
2         def __init__(self, window_size):
3             self.window_size = window_size
4             self.window = []
5
6         def add_data(self, data):
7             self.window.append(data)
8             # Slide of one position if the window
9             # is full
10            if len(self.window) >
11                self.window_size:
12                self.window.pop(0)
13
14        def process_window(self):
15            # Perform processing on the window
16            # data
17            # Example: Calculate the average of
18            # the window
19            window_sum = sum(self.window)
20            window_avg = window_sum /
21                len(self.window)
22            return window_avg
23
24    # Example usage
25    window_size = 5
26    stream = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
27    sliding_window = SlidingWindow(window_size)

```

```
23
24     for data in stream:
25         sliding_window.add_data(data)
26         window_avg =
27             sliding_window.process_window()
28         print(f"Window Average: {window_avg}")
```

Naturally, as easily verified, the output in this case will be:

```
1      Window Average: 1.0
2      Window Average: 1.5
3      Window Average: 2.0
4      Window Average: 2.5
5      Window Average: 3.0
6      Window Average: 4.0
7      Window Average: 5.0
8      Window Average: 6.0
9      Window Average: 7.0
10     Window Average: 8.0
```