

Bucket Sort

Floris Videler

1758374

5-2-2020

Inleiding

Voor het vak high performance programming was de eerste opdracht het maken van een variant op het bucket sort sorteer algoritme en deze analyseren.

Algoritme

Mijn code volgt de pseudo gegeven op canvas. Mijn implementatie verdeelt de lijst eerst onder in twee lijsten, eentje voor de positieve getallen en eentje voor de negatieve getallen. Deze worden daarna apart in mijn recursieve bucket sort functie gegooit.

Analyse

Voor we beginnen met de analyse is het handig om stil te staan bij twee letters en wat deze betekenen:

1. n : Het aantal elementen in de input lijst.
2. k : De lengte van het langste item in de lijst.

In het bestand "*analyze.py*" kan een programmatische analyse gevonden worden die de onderstaande resultaten onderbouwd.

Tijdscomplexiteit

Het analyseren van zelf geschreven code was een stuk lastiger dan gedacht. Hier onder staan twee afbeeldingen, de eerste is de functie is ook wel de "verdeel" functie. De tweede afbeelding is van het algoritme zelf.

```
32 def bucket_sort(data: list) -> list:
33     """
34     Splits a list in negative and positive numbers and bucket sorts them individually.
35     :param data: List of unsorted ints.
36     :return: Sorted list
37     """
38     pos = [] 1. Variable declaration: 1
39     neg = [] 2. Variable declaration: 1
40     # Separate the positive numbers from the negative numbers.
41     for i in data: 3. Loop over n item: n
42         if i < 0: 4. If statement: 1
43             neg.append(i)
44         else: 5. Append statement: 1
45             pos.append(i)
46     # Before returning the combination of the 2 lists we reverse the negative list.
47     # We also subtract 1 from length of the longest int because of the "-" symbol.
48     return list(reversed(bucket_sort_function(neg, len(str(min(neg))) - 1)) + bucket_sort_function(pos,
49     6. Min & max function: 1/2n + 1/2n = n 7. Str function: 2*k (2* str functie) len(str(max(pos)))
50     8. Len function: 1*2=2
51     24. Reverse function: n 25. List function: n
52     3n
53
54 def bucket_sort_function(data: list, k: int = None, p: int = 0) -> list:
55     """
56     This function sorts a list using bucket sort. Used by the bucket_sort function.
57     :param data: The List of ints to sort.
58     :param k: The Length of the Largest number in the List.
59     :param p: Times the function has run.
60     :return: A sorted List.
61     """
62     # The stop condition, if the times we run the sort is equal to the max size of the ints we stop.
63     if p * -1 == k: 9. Times operation: 1 10. Vergelijking: 1
64         return data
65     # Init the bucket list
66     bucket_list = [[]] 11. Variable declaration: 1
67     for i in data[:]: 12. For loop over n items: n
68         # We try to get the right int to check what bucket is needed. Is this fails for some reason
69         # (example: "-" is found instead of an int) we do leave the data where it is.
70         try: 13. Try except: 1
71             check = int(str(i)[p - 1]) 14. Variable declaration: 1 15. Str function: k 16. Get by index: 1 17. Int function: k
72             bucket_list[check].append(i) 18. Get by index: 1 19. Append: 1
73             data.remove(i) 20. Remove: n2
74         except:
75             pass
76     # Add data from buckets back to the original data.
77     for i in bucket_list: 21. Loop for 10 items: 10
78         data += i 22. Extend list: 1/10 n (De items zijn gelijk verdeeld over de 10 buckets) } 10 * 1/10n = n
79     # Recursively repeat this. 23. Repeat k times
80     return bucket_sort_function(data, k, p - 1)
```

Hier onder heb ik complexiteit van iedere stap uitgeschreven:

$$1 + 1 + 3n + n + 2k + 2 + 1 + k(1 + 1 + n(1 + 1 + k + 1 + k + 1 + 1 + n) + n) + n + n$$

Vereenvoudigd:

$$1 + 1 + 6n + 2k + 2 + 1 + k(1 + 1 + 2n(1 + 1 + 2k + 1 + 1 + 1 + n))$$

Weghalen van alle constanten:

$$n + k + k (n (k + n))$$

Haakjes weghalen:

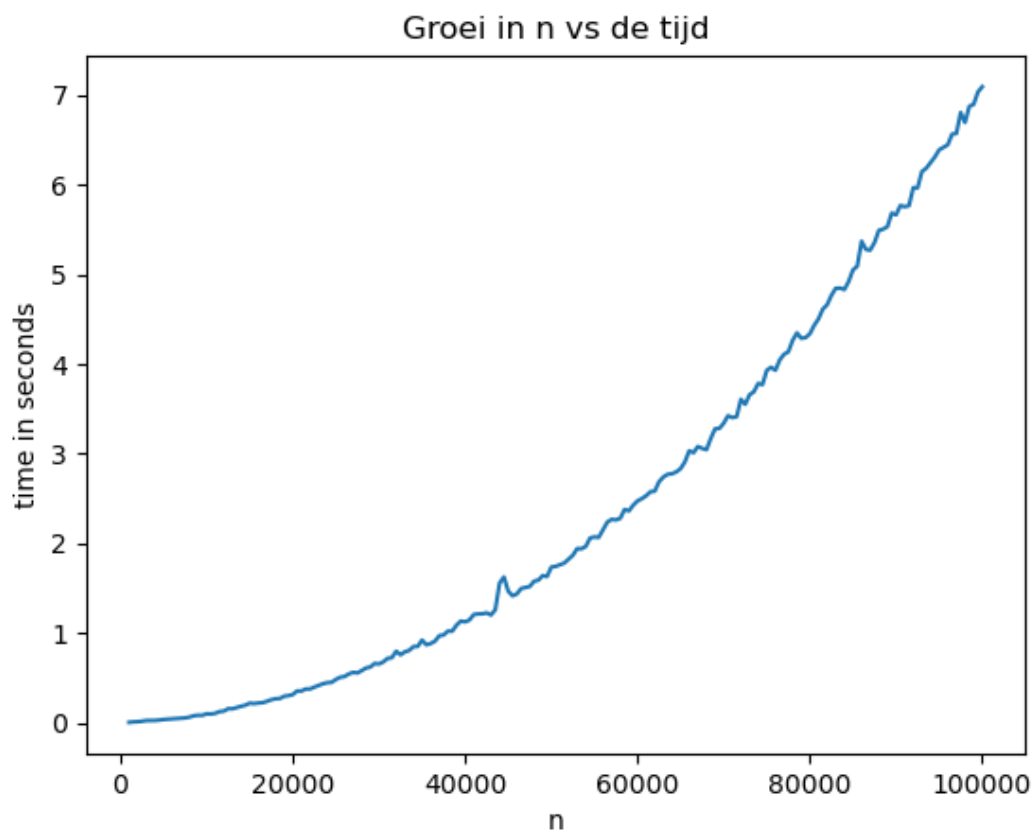
$$n + k + k (n (k + n))$$

$$n + k + k (nk + n^2)$$

$$n + k + k^2n + kn^2$$

Als we hier de dominante term uit willen halen lopen we tegen een probleem aan. k^2n kan groter zijn dan kn^2 afhankelijk van de grote van k en n . Deze kunnen beide groeien tot oneindig. Als er ik toch één moet kiezen, zou ik kiezen voor kn^2 . Dit omdat n vaak relevanter bij het bepalen van tijdcomplexiteit en de big O notatie. De big O notatie zou in dit geval dan zijn: $O(kn^2)$.

De onderstaande grafiek bevestigt ook dat de complexiteit inderdaad wel eens exponentieel zou kunnen zijn.



Wat natuurlijk ook effect kan hebben op de snelheid van een algoritme is de input: zijn er dubbele getallen? Is de lijst al gesorteerd? Etc. In het bestand “analyze.py” heb ik hier onderzoek naar gedaan. In de onderstaande afbeelding is hier een korte samenvatting van te zien.

Each metric is calculated with a population of 100 and a list length of 75000

Type	Avg	Min	Max	Std
-----	-----	-----	-----	-----
Bucket sort normal	3.18761	3.02368	3.50576	0.0574777
Bucket sort sorted list	3.18882	3.08348	3.39032	0.0452626
bucket sort reversed list	3.18486	3.04703	3.31774	0.0472281
bucket sort unique values	3.00869	2.91375	3.22659	0.0451142

In de tabel wordt voor ieder geval het volgende getoond:

- Gemiddelde
- Kortste tijd
- Langste tijd
- En de standaardafwijking

Op basis van deze gegevens kunnen we de conclusie trekken dat geen speciaal geval echt effect heeft op de performance van het algoritme. Het verschil dat gezien wordt is zeker geen groot verschil en kan ook toeval worden genoemd.

Ruimtecomplexiteit

Over het algemeen wordt er gezegd dat ruimtecomplexiteit minder belangrijk is dan tijd complexiteit. We hebben hier ook een stuk minder uitleg over gekregen. In mijn algoritme heb ik een ruimtecomplexiteit bepaald van $O(kn)$. Dit omdat in het slechtste geval alle elementen in de lijst k keer worden ingedeeld in buckets. Dit zou betekenen dat er k keer nieuwe buckets volledig gevuld zouden moeten worden met n items.

Conclusie en discussie

Terug kijkend om mijn uitwerking en analyse hiervan zijn er een aantal punten die genoemd moeten worden:

1. Mijn uitwerking is niet de beste: Ik heb mijn beste gedaan om dit algoritme zo efficiënt mogelijk te maken en ongetwijfeld is er ook een efficiëntere manier. Dit kan dus leiden tot een onjuiste analyse voor bucket sort. Het is dus belangrijk te onthouden dat deze analyse is gemaakt voor mijn versie van het algoritme.
2. Er is zoveel mogelijk meegenomen in het bepalen van de big O: In de les was er al besproken wat er wel en niet meegenomen moet worden in het bepalen van de big O. Tellen variable assignments mee? Tellen normale Python functies mee? Hier is niet één definitief antwoord op. Voor mijn analyse heb ik dit wel allemaal meegenomen. Dit heb ik gedaan om de analyse zo precies mogelijk te laten kloppen met de code die er staat. Een goed voorbeeld hier van is de remove functie voor een Python lijst, deze heeft een complexiteit van $O(n)$. Deze functie in een for loop gebruiken zou de complexiteit al snel verhogen.
3. Hoe bruikbaar is deze implementatie? Mijn specifieke implementatie van deze variatie op bucket sort is zeker niet het beste sorteer algoritme. Er zijn veel algoritmes met een betere tijd en ruimtecomplexiteit. Wel vind ik het belangrijk

om te noemen dat ondanks mijn uitwerking niet de snelste of de beste is, ik wel anders ben gaan kijken naar hoe je algoritmes bouwt en welke aanpak je gebruikt. Het opzoeken van complexiteiten van Python functies heeft mij veel kennis gegeven over wat voor effect input heeft op de uitvoer tijd.