#### Analysis of algorithms

#### Madhavan Mukund

#### https://www.cmi.ac.in/~madhavan

# Programming and Data Structures with Python Lecture 15, 6 Oct 2022

- Every SIM card needs to be linked to an Aadhaar card
- Validate Aadhaar number for each SIM card

- Every SIM card needs to be linked to an Aadhaar card
- Validate Aadhaar number for each SIM card
- Simple nested loop

for each SIM card S:
for each Aadhaar number A:
 check if Aadhaar number in S
 matches A

- Every SIM card needs to be linked to an Aadhaar card
- Validate Aadhaar number for each SIM card
- Simple nested loop
- How long will this take?
  - M SIM cards, N Aadhaar cards
  - Nested loops iterate M · N times

for each SIM card S: M for each Aadhaar number A: N check if Aadhaar number in S matches A

- Every SIM card needs to be linked to an Aadhaar card
- Validate Aadhaar number for each SIM card
- Simple nested loop
- How long will this take?
  - M SIM cards, N Aadhaar cards
  - Nested loops iterate M · N times
- What are M and N
  - Almost everyone in India has an Aadhaar card: N > 10<sup>9</sup>
  - Number of SIM cards registered is similar: M > 10<sup>9</sup>

```
for each SIM card S:
for each Aadhaar number A:
  check if Aadhaar number in S
  matches A
```

- Assume  $M = N = 10^9$
- Nested loops execute 10<sup>18</sup> times

for each SIM card S:
for each Aadhaar number A:
 check if Aadhaar number in S
 matches A

- Assume  $M = N = 10^9$
- Nested loops execute 10<sup>18</sup> times
- We calculated that Python can perform 10<sup>7</sup> operations in a second
- This will take at least 10<sup>11</sup> seconds

```
for each SIM card S:
for each Aadhaar number A:
  check if Aadhaar number in S
  matches A
```

- Assume  $M = N = 10^9$
- Nested loops execute 10<sup>18</sup> times
- We calculated that Python can perform 10<sup>7</sup> operations in a second
- This will take at least 10<sup>11</sup> seconds
   10<sup>11</sup>/60 ≈ 1.67 × 10<sup>9</sup> minutes

for each SIM card S:
for each Aadhaar number A:
 check if Aadhaar number in S
 matches A

- Assume  $M = N = 10^9$
- Nested loops execute 10<sup>18</sup> times
- We calculated that Python can perform 10<sup>7</sup> operations in a second
- This will take at least 10<sup>11</sup> seconds
  - $10^{11}/60 \approx 1.67 \times 10^9$  minutes
  - $\blacksquare~(1.67\times 10^9)/60\approx 2.8\times 10^7~hours$

```
for each SIM card S:
for each Aadhaar number A:
  check if Aadhaar number in S
  matches A
```

- Assume  $M = N = 10^9$
- Nested loops execute 10<sup>18</sup> times
- We calculated that Python can perform 10<sup>7</sup> operations in a second
- This will take at least 10<sup>11</sup> seconds
  - $10^{11}/60 \approx 1.67 \times 10^9$  minutes
  - $(1.67 \times 10^9)/60 \approx 2.8 \times 10^7$  hours
  - $\blacksquare~(2.8\times10^7)/24\approx1.17\times10^6~days$

```
for each SIM card S:
for each Aadhaar number A:
  check if Aadhaar number in S
  matches A
```

PDSP Lecture 15 3/1

- Assume  $M = N = 10^9$
- Nested loops execute 10<sup>18</sup> times
- We calculated that Python can perform 10<sup>7</sup> operations in a second
- This will take at least 10<sup>11</sup> seconds
  - $10^{11}/60 \approx 1.67 \times 10^9$  minutes
  - $(1.67 \times 10^9)/60 \approx 2.8 \times 10^7$  hours
  - $\blacksquare~(2.8\times10^7)/24\approx1.17\times10^6~days$
  - $(1.17 \times 10^6)/365 \approx 3200$  years!

```
for each SIM card S:
for each Aadhaar number A:
  check if Aadhaar number in S
  matches A
```

- Assume  $M = N = 10^9$
- Nested loops execute 10<sup>18</sup> times
- We calculated that Python can perform 10<sup>7</sup> operations in a second
- This will take at least 10<sup>11</sup> seconds
  - $10^{11}/60 \approx 1.67 \times 10^9$  minutes
  - $(1.67\times 10^9)/60\approx 2.8\times 10^7$  hours
  - $(2.8 \times 10^7)/24 \approx 1.17 \times 10^6$  days
  - $(1.17 \times 10^6)/365 \approx 3200$  years!
- How can we fix this?

```
for each SIM card S:
for each Aadhaar number A:
  check if Aadhaar number in S
  matches A
```



- You propose a date
- I answer, Yes, Earlier, Later
- Suppose my birthday is 12 April

- You propose a date
- I answer, Yes, Earlier, Later
- Suppose my birthday is 12 April
- A possible sequence of questions
  - September 12? Earlier
  - February 23? Later
  - July 2? Earlier
  - ...

- You propose a date
- I answer, Yes, Earlier, Later
- Suppose my birthday is 12 April
- A possible sequence of questions
  - September 12? Earlier
  - February 23? Later
  - July 2? Earlier
  - • •
- What is the best strategy?

- You propose a date
- I answer, Yes, Earlier, Later
- Suppose my birthday is 12 April
- A possible sequence of questions
  - September 12? Earlier
  - February 23? Later
  - July 2? Earlier
  - • •
- What is the best strategy?

#### Interval of possibilities

- You propose a date
- I answer, Yes, Earlier, Later
- Suppose my birthday is 12 April
- A possible sequence of questions
  - September 12? Earlier
  - February 23? Later
  - July 2? Earlier
  - • •
- What is the best strategy?

- Interval of possibilities
- Query midpoint halves the interval
  - June 30? Earlier
  - March 31? Later
  - May 15? Earlier
  - April 22? Earlier
  - April 11? Later
  - April 16? Earlier
  - April 13? Earlier
  - April 12? Yes

- You propose a date
- I answer, Yes, Earlier, Later
- Suppose my birthday is 12 April
- A possible sequence of questions
  - September 12? Earlier
  - February 23? Later
  - July 2? Earlier
  - • •
- What is the best strategy?

- Interval of possibilities
- Query midpoint halves the interval
  - June 30? Earlier
  - March 31? Later
  - May 15? Earlier
  - April 22? Earlier
  - April 11? Later
  - April 16? Earlier
  - April 13? Earlier
  - April 12? Yes
- Under 10 questions

- Assume Aadhaar details are sorted by Aadhaar number
- Use the halving strategy to check each SIM card

for each SIM card S:
 probe sorted Aadhaar list to
 find a match with S

- Assume Aadhaar details are sorted by Aadhaar number
- Use the halving strategy to check each SIM card
- Halving 10 times reduces the interval by a factor of 1000, because 2<sup>10</sup> = 1024

```
for each SIM card S:
  probe sorted Aadhaar list to
  find a match with S
```

- Assume Aadhaar details are sorted by Aadhaar number
- Use the halving strategy to check each SIM card
- Halving 10 times reduces the interval by a factor of 1000, because 2<sup>10</sup> = 1024
- After 10 queries, interval shrinks to 10<sup>6</sup>
- After 20 queries, interval shrinks to  $10^3$
- After 30 queries, interval shrinks to 1

```
for each SIM card S:
  probe sorted Aadhaar list to
  find a match with S
```

 $2^{10} = 10^{3}$   $10^{7} = 10^{3} \times 10^{3} \times 10^{3}$   $= 2^{10} \times 2^{10} \times 2^{10}$   $= 2^{30}$ 

- Assume Aadhaar details are sorted by Aadhaar number
- Use the halving strategy to check each SIM card
- Halving 10 times reduces the interval by a factor of 1000, because 2<sup>10</sup> = 1024
- After 10 queries, interval shrinks to 10<sup>6</sup>
- After 20 queries, interval shrinks to 10<sup>3</sup>
- After 30 queries, interval shrinks to 1
- Total time  $\approx 10^9 \times 30$

```
for each SIM card S:
  probe sorted Aadhaar list to
  find a match with S
```

PDSP Lecture 15 5/1

- Assume Aadhaar details are sorted by Aadhaar number
- Use the halving strategy to check each SIM card
- Halving 10 times reduces the interval by a factor of 1000, because 2<sup>10</sup> = 1024
- After 10 queries, interval shrinks to 10<sup>6</sup>
- After 20 queries, interval shrinks to 10<sup>3</sup>
- After 30 queries, interval shrinks to 1
- Total time  $\approx 10^9 \times 30$

```
for each SIM card S:
  probe sorted Aadhaar list to
  find a match with S
```

- 3000 seconds, or 50 minutes
- From 3200 years to 50 minutes!

- Assume Aadhaar details are sorted by Aadhaar number
- Use the halving strategy to check each SIM card
- Halving 10 times reduces the interval by a factor of 1000, because 2<sup>10</sup> = 1024
- After 10 queries, interval shrinks to 10<sup>6</sup>
- After 20 queries, interval shrinks to 10<sup>3</sup>
- After 30 queries, interval shrinks to 1
- Total time  $\approx 10^9 \times 30$

```
for each SIM card S:
  probe sorted Aadhaar list to
  find a match with S
```

- 3000 seconds, or 50 minutes
- From 3200 years to 50 minutes!
- Of course, to achieve this we have to first sort the Aadhaar cards
- Arranging the data results in a much more efficient solution
- Both algorithms and data structures matter

#### Comparing orders of magnitude

Madhavan Mukund

https://www.cmi.ac.in/~madhavan

Programming and Data Structures with Python Lecture 15, 6 Oct 2022

・ロト ・日ト ・ヨト ・ヨト ・ヨー うへで

# Orders of magnitude

- When comparing t(n), focus on orders of magnitude
  - Ignore constant factors
- $f(n) = n^3$  eventually grows faster than  $g(n) = 5000n^2$

# Orders of magnitude

- When comparing t(n), focus on orders of magnitude
  - Ignore constant factors
- $f(n) = n^3$  eventually grows faster than  $g(n) = 5000n^2$
- How do we compare functions with respect to orders of magnitude?

### Upper bounds

 f(x) is said to be O(g(x)) if we can find constants c and x₀ such that c ⋅ g(x) is an upper bound for f(x) for x beyond x₀

- 4 西

э

### Upper bounds

- f(x) is said to be O(g(x)) if we can find constants c and x₀ such that c ⋅ g(x) is an upper bound for f(x) for x beyond x₀
- $f(x) \leq cg(x)$  for every  $x \geq x_0$



## Upper bounds

- f(x) is said to be O(g(x)) if we can find constants c and x₀ such that c ⋅ g(x) is an upper bound for f(x) for x beyond x₀
- $f(x) \leq cg(x)$  for every  $x \geq x_0$
- Graphs of typical functions we have seen



## Examples

#### ■ 100n + 5 is $O(n^2)$

- $100n + 5 \le 100n + n = 101n$ , for  $n \ge 5$
- $101n \le 101n^2$
- Choose  $n_0 = 5$ , c = 101



#### Examples

■ 100n + 5 is  $O(n^2)$ 

- $100n + 5 \le 100n + n = 101n$ , for  $n \ge 5$
- $101n \le 101n^2$
- Choose  $n_0 = 5$ , c = 101
- Alternatively
  - $100n + 5 \le 100n + 5n = 105n$ , for  $n \ge 1$
  - $105n \le 105n^2$
  - Choose  $n_0 = 1$ , c = 105



#### Examples

■ 100n + 5 is  $O(n^2)$ 

- $100n + 5 \le 100n + n = 101n$ , for  $n \ge 5$
- $101n \le 101n^2$
- Choose  $n_0 = 5$ , c = 101
- Alternatively
  - $100n + 5 \le 100n + 5n = 105n$ , for  $n \ge 1$
  - $105n \le 105n^2$
  - Choose  $n_0 = 1$ , c = 105
- Choice of  $n_0$ , c not unique



#### Examples . . .

- $100n^2 + 20n + 5$  is  $O(n^2)$ 
  - $100n^2 + 20n + 5 \le 100n^2 + 20n^2 + 5n^2$ , for  $n \ge 1$
  - $100n^2 + 20n + 5 \le 125n^2$ , for  $n \ge 1$
  - Choose  $n_0 = 1$ , c = 125



- $100n^2 + 20n + 5$  is  $O(n^2)$ 
  - $100n^2 + 20n + 5 \le 100n^2 + 20n^2 + 5n^2$ , for  $n \ge 1$
  - $100n^2 + 20n + 5 \le 125n^2$ , for  $n \ge 1$
  - Choose  $n_0 = 1$ , c = 125
- What matters is the highest term
  - 20n + 5 is dominated by  $100n^2$



■  $100n^2 + 20n + 5$  is  $O(n^2)$ 

- $100n^2 + 20n + 5 \le 100n^2 + 20n^2 + 5n^2$ , for  $n \ge 1$
- $100n^2 + 20n + 5 \le 125n^2$ , for  $n \ge 1$
- Choose  $n_0 = 1$ , c = 125
- What matters is the highest term
  - 20n + 5 is dominated by  $100n^2$
- $n^3$  is not  $O(n^2)$ 
  - No matter what c we choose, cn<sup>2</sup> will be dominated by n<sup>3</sup> for n ≥ c



If f<sub>1</sub>(n) is O(g<sub>1</sub>(n)) and f<sub>2</sub>(n) is O(g<sub>2</sub>(n)), then f<sub>1</sub>(n) + f<sub>2</sub>(n) is O(max(g<sub>1</sub>(n), g<sub>2</sub>(n)))

m + n

イロト 不得下 イヨト イヨト

<u>m > n</u>



1 つ 1

 $\leq 2n$ 

3

• If  $f_1(n)$  is  $O(g_1(n))$  and  $f_2(n)$  is  $O(g_2(n))$ , then  $f_1(n) + f_2(n)$  is  $O(\max(g_1(n), g_2(n)))$ 

Proof

- $f_1(n) \le c_1 g_1(n)$  for  $n > n_1$
- $f_2(n) \le c_2 g_2(n)$  for  $n > n_2$

- 4 西

э

• If  $f_1(n)$  is  $O(g_1(n))$  and  $f_2(n)$  is  $O(g_2(n))$ , then  $f_1(n) + f_2(n)$  is  $O(\max(g_1(n), g_2(n)))$ 

- $f_1(n) \le c_1 g_1(n)$  for  $n > n_1$
- $f_2(n) \le c_2 g_2(n)$  for  $n > n_2$
- Let  $c_3 = \max(c_1, c_2)$ ,  $n_3 = \max(n_1, n_2)$

• If  $f_1(n)$  is  $O(g_1(n))$  and  $f_2(n)$  is  $O(g_2(n))$ , then  $f_1(n) + f_2(n)$  is  $O(\max(g_1(n), g_2(n)))$ 

- $f_1(n) \le c_1 g_1(n)$  for  $n > n_1$
- $f_2(n) \le c_2 g_2(n)$  for  $n > n_2$
- Let  $c_3 = \max(c_1, c_2)$ ,  $n_3 = \max(n_1, n_2)$
- For  $n \ge n_3$ ,  $f_1(n) + f_2(n)$ 
  - $\leq c_1g_1(n)+c_2g_2(n)$

• If  $f_1(n)$  is  $O(g_1(n))$  and  $f_2(n)$  is  $O(g_2(n))$ , then  $f_1(n) + f_2(n)$  is  $O(\max(g_1(n), g_2(n)))$ 

- $f_1(n) \le c_1 g_1(n)$  for  $n > n_1$
- $f_2(n) \le c_2 g_2(n)$  for  $n > n_2$
- Let  $c_3 = \max(c_1, c_2)$ ,  $n_3 = \max(n_1, n_2)$
- For  $n \ge n_3$ ,  $f_1(n) + f_2(n)$ 
  - $\leq c_1g_1(n)+c_2g_2(n)$
  - $\leq c_3(g_1(n)+g_2(n))$

• If  $f_1(n)$  is  $O(g_1(n))$  and  $f_2(n)$  is  $O(g_2(n))$ , then  $f_1(n) + f_2(n)$  is  $O(\max(g_1(n), g_2(n)))$ 

- $f_1(n) \le c_1 g_1(n)$  for  $n > n_1$
- $f_2(n) \le c_2 g_2(n)$  for  $n > n_2$
- Let  $c_3 = \max(c_1, c_2)$ ,  $n_3 = \max(n_1, n_2)$
- For  $n \ge n_3$ ,  $f_1(n) + f_2(n)$ 
  - $\leq c_1g_1(n)+c_2g_2(n)$
  - $\leq c_3(g_1(n)+g_2(n))$
  - $\leq 2c_3(\max(g_1(n),g_2(n)))$

• If  $f_1(n)$  is  $O(g_1(n))$  and  $f_2(n)$  is  $O(g_2(n))$ , then  $f_1(n) + f_2(n)$  is  $O(\max(g_1(n), g_2(n)))$ 

- $f_1(n) \le c_1 g_1(n)$  for  $n > n_1$
- $f_2(n) \le c_2 g_2(n)$  for  $n > n_2$
- Let  $c_3 = \max(c_1, c_2)$ ,  $n_3 = \max(n_1, n_2)$
- For  $n \ge n_3$ ,  $f_1(n) + f_2(n)$ 
  - $\leq c_1g_1(n)+c_2g_2(n)$
  - $\leq c_3(g_1(n)+g_2(n))$
  - $\leq 2c_3(\max(g_1(n),g_2(n)))$

- Algorithm has two phases
  - Phase A takes time  $O(g_A(n))$
  - Phase B takes time  $O(g_B(n))$

• If  $f_1(n)$  is  $O(g_1(n))$  and  $f_2(n)$  is  $O(g_2(n))$ , then  $f_1(n) + f_2(n)$  is  $O(\max(g_1(n), g_2(n)))$ 

- $f_1(n) \le c_1 g_1(n)$  for  $n > n_1$
- $f_2(n) \le c_2 g_2(n)$  for  $n > n_2$
- Let  $c_3 = \max(c_1, c_2)$ ,  $n_3 = \max(n_1, n_2)$
- For  $n \ge n_3$ ,  $f_1(n) + f_2(n)$ 
  - $\leq c_1g_1(n)+c_2g_2(n)$
  - $\leq c_3(g_1(n)+g_2(n))$
  - $\leq 2c_3(\max(g_1(n),g_2(n)))$

- Algorithm has two phases
  - Phase A takes time  $O(g_A(n))$
  - Phase B takes time  $O(g_B(n))$
- Algorithm as a whole takes time max(O(g<sub>A</sub>(n), g<sub>B</sub>(n)))

• If  $f_1(n)$  is  $O(g_1(n))$  and  $f_2(n)$  is  $O(g_2(n))$ , then  $f_1(n) + f_2(n)$  is  $O(\max(g_1(n), g_2(n)))$ 

- $f_1(n) \le c_1 g_1(n)$  for  $n > n_1$
- $f_2(n) \le c_2 g_2(n)$  for  $n > n_2$
- Let  $c_3 = \max(c_1, c_2)$ ,  $n_3 = \max(n_1, n_2)$
- For  $n \ge n_3$ ,  $f_1(n) + f_2(n)$ 
  - $\leq c_1g_1(n)+c_2g_2(n)$
  - $\leq c_3(g_1(n)+g_2(n))$
  - $\leq 2c_3(\max(g_1(n),g_2(n)))$

- Algorithm has two phases
  - Phase A takes time  $O(g_A(n))$
  - Phase B takes time  $O(g_B(n))$
- Algorithm as a whole takes time max(O(g<sub>A</sub>(n), g<sub>B</sub>(n)))
- Least efficient phase is the upper bound for the whole algorithm

#### Lower bounds

- f(x) is said to be Ω(g(x)) if we can find constants c and x<sub>0</sub> such that cg(x) is a lower bound for f(x) for x beyond x<sub>0</sub>
  - $f(x) \ge cg(x)$  for every  $x \ge x_0$

- 4 西

э

#### Lower bounds

- f(x) is said to be Ω(g(x)) if we can find constants c and x<sub>0</sub> such that cg(x) is a lower bound for f(x) for x beyond x<sub>0</sub>
  - $f(x) \ge cg(x)$  for every  $x \ge x_0$
- $n^3$  is  $\Omega(n^2)$ 
  - $n^3 > n^2$  for all *n*, so  $n_0 = 1$ , c = 1

#### Lower bounds

- f(x) is said to be Ω(g(x)) if we can find constants c and x<sub>0</sub> such that cg(x) is a lower bound for f(x) for x beyond x<sub>0</sub>
  - $f(x) \ge cg(x)$  for every  $x \ge x_0$
- $n^3$  is  $\Omega(n^2)$ 
  - $n^3 > n^2$  for all *n*, so  $n_0 = 1$ , c = 1
- Typically we establish lower bounds for a problem rather than an individual algorithm
  - If we sort a list by comparing elements and swapping them, we require  $\Omega(n \log n)$  comparisons
  - This is independent of the algorithm we use for sorting

- f(x) is said to be  $\Theta(g(x))$  if it is both O(g(x)) and  $\Omega(g(x))$ 
  - Find constants  $c_1, c_2, x_0$  such that  $c_1g(x) \le f(x) \le c_2g(x)$  for every  $x \ge x_0$

э

Image: A math a math

- f(x) is said to be  $\Theta(g(x))$  if it is both O(g(x)) and  $\Omega(g(x))$ 
  - Find constants  $c_1, c_2, x_0$  such that  $c_1g(x) \le f(x) \le c_2g(x)$  for every  $x \ge x_0$

■ n(n-1)/2 is  $\Theta(n^2)$ 

< /□ > < 三

3

- f(x) is said to be  $\Theta(g(x))$  if it is both O(g(x)) and  $\Omega(g(x))$ 
  - Find constants  $c_1, c_2, x_0$  such that  $c_1g(x) \le f(x) \le c_2g(x)$  for every  $x \ge x_0$

■ n(n-1)/2 is  $\Theta(n^2)$ 

Upper bound

■ 
$$n(n-1)/2 = n^2/2 - n/2 \le n^2/2$$
 for all  $n \ge 0$ 

э

- f(x) is said to be  $\Theta(g(x))$  if it is both O(g(x)) and  $\Omega(g(x))$ 
  - Find constants  $c_1, c_2, x_0$  such that  $c_1g(x) \le f(x) \le c_2g(x)$  for every  $x \ge x_0$

■ n(n-1)/2 is  $\Theta(n^2)$ 

- Upper bound
  - $n(n-1)/2 = n^2/2 n/2 \le n^2/2$  for all  $n \ge 0$
- Lower bound

■ 
$$n(n-1)/2 = n^2/2 - n/2 \ge n^2/2 - (n/2 \times n/2) \ge n^2/4$$
 for  $n \ge 2$ 

э

- f(x) is said to be  $\Theta(g(x))$  if it is both O(g(x)) and  $\Omega(g(x))$ 
  - Find constants  $c_1, c_2, x_0$  such that  $c_1g(x) \le f(x) \le c_2g(x)$  for every  $x \ge x_0$

■ n(n-1)/2 is  $\Theta(n^2)$ 

- Upper bound
  - $n(n-1)/2 = n^2/2 n/2 \le n^2/2$  for all  $n \ge 0$
- Lower bound

■ 
$$n(n-1)/2 = n^2/2 - n/2 \ge n^2/2 - (n/2 \times n/2) \ge n^2/4$$
 for  $n \ge 2$ 

• Choose  $n_0 = 2$ ,  $c_1 = 1/4$ ,  $c_2 = 1/2$ 

3



• f(n) is O(g(n)) means g(n) is an upper bound for f(n)

Useful to describe asymptotic worst case running time



What is input size? Multiply a × b Cost grows with values of a & b  $2^{\chi} = 10^{\gamma}$ 3578 N. J drzik  $\times \log_2 2 = y \log_2 10$ x 764\$\$9 logion x = y - log to rdd L \_ x Dgze rdd L \_ x  $= C - \log_2 n$ 

# Summary

- f(n) is O(g(n)) means g(n) is an upper bound for f(n)
  - Useful to describe asymptotic worst case running time
- f(n) is  $\Omega(g(n))$  means g(n) is a lower bound for f(n)
  - Typically used for a problem as a whole, rather than an individual algorihm

# Summary

- f(n) is O(g(n)) means g(n) is an upper bound for f(n)
  - Useful to describe asymptotic worst case running time
- f(n) is  $\Omega(g(n))$  means g(n) is a lower bound for f(n)
  - Typically used for a problem as a whole, rather than an individual algorihm
- f(n) is  $\Theta(g(n))$ : matching upper and lower bounds
  - We have found an optimal algorithm for a problem