COMPUTATIONAL EFFICIENCY

LING 1340/2340: Data Science for Linguists

Joey Livorno

Objectives

- Computer Anatomy
 - Memory
 - Processing
- Computational Efficiency
 - Algorithmic complexity
 - Big O notation
- Solutions
 - Chunking data
 - Choosing the right data type

How Do Computers Work?

- Main Components
 - CPU
 - Video Card
 - RAM (Random Access Memory)
 - Hard Drive
 - Motherboard
 - Power Supply



Random Access Memory (RAM)

- "Random Access"
 - Can be read or changed in any order
 - Makes this type of storage insanely fast
 - Used for our "working" data
- Other storage mediums
 - HDD, SSD, Magnetic Tape, etc.
 - Not as fast
 - Used for large, permanent storage



Central Processing Unit (CPU)

- Circuit that executes instructions of programs
 - Arithmetic
 - Logic
 - Controlling
 - I/O
- Contain Cores
 - Individual Processors within the CPU
 - Each core can compute independently
 - Allows the CPU to process multiple tasks in parallel





Computational Efficiency: Space vs. Time

Space Complexity

- How much memory are we using?
- Do not create duplicate data objects
- Avoid creating data objects that do not need to be stored in their entirety
- Garbage Collection
 - Variable name acts as pointer
 - Data that does not have a pointer will be removed from memory automatically
 - If space is a concern, don't rely on this

Time Complexity

- How much processing are we doing?
- Avoid repeating steps that require heavy processing
- Use efficient algorithms
- Use data types optimal for a specific task

The goal is to optimize both. They often exist as a trade-off relationship, but by managing our resources we can achieve balance.

Algorithmic Complexity and the Big O

- <u>https://rob-bell.net/2009/06/a-beginners-guide-to-big-o-notation/</u>
- Can be used to describe the worst-case performance of an algorithm
- We have an **unsorted list of** *n* **items**. Imagine *n* to be any value.
- 1. Is the first element in the list an even number?
 - Can be implemented in O(1), or constant time
 - Does not change as the size of n increases
- 2. Does the list contain 42?
 - Can be implemented in O(n), or linear time
 - Directly proportional to the size of n
- 3. Does the list contain duplicate values?
 - Can be implemented in O(n²), or quadratic time
 - Directly proportional to the square of the size of n
 - Can be optimized, but at the expense of space complexity

Algorithmic Complexity and the Big O

- 4. Sort the list
 - Can be implemented in O(n log n), or loglinear time
 - See: <u>https://brilliant.org/wiki/sorting-algorithms/</u>



Opening and Closing Files

```
f = open('review.json')
lines = f.readlines()
for l in lines:
    if 'horrible' in l:
        print(l)
f.close()
```

```
f = open('review.json')
for l in f:
    if 'horrible' in l:
        print(l)
f.close()
```

lines = open('review.json').readlines()
for l in lines :
 if 'horrible' in l:
 print(l)

with open('review.json') as f: for l in f: if 'horrible' in l: print(l)

Which methods are memoryefficient? Answer: the ones on the right are more memory-efficient.

Left two boxes create lines object which holds the entirety of the .json file in memory!

Processing Big Files

• How much memory did we use for process_reviews.py?

import pandas as pd	
import sys from collections import Counter	df holds the entire ison
<pre>filename = sys.argv[1] # 5 GBs</pre>	file's content as DF, in memory.
df = pd.read_json(filename, lines=True, encoding='utf-8') print(df.head(5)) # ~4 GBs	wtoks is another big data object created from the entire text.
<pre>wtoks = ' '.join(df['text']).split() wfreq = Counter(wtoks) print(wfreq.most_common(20)) ^G Get Help ^O WriteOut ^R Read File ^Y Prev Page ^K Cut Text ^C Cur Pos</pre>	Processing 4-million Yelp reviews took 6.5 minutes and 33.5 GB of memory.

Handling Files in Chunks

- We can store the data in chunks so that we only have access to what we need at a given point in time
- <u>https://realpython.com/introduction-to-python-generators/</u>

	🚸 naraehan@login0:~/todo11-12 — 🗌	×
	GNU nano 2.3.1 File: process_reviews_eff.py	<u>^</u>
Reading in the files as a generator takes up zero space.	<pre>import pandas as pd import sys from collections import Counter filename = sys.argv[1] df_chunks = pd.read_json(filename, chunksize=10000, lines=True, encoding='utf-8</pre>	Read_json() argument says that each chunk will be 10k lines.
	wfreq = Counter()	
Then we simply iterate through the smaller dfs.	for chunk in df_chunks: for text in chunk['text']: wfreq.update(text.split())	This is memory efficient!
	<pre>print(wfreq.most_common(20))</pre>	Processing the same 4million Yelp reviews
	∧G Get Help ∧O WriteOut ∧R Read File ∧Y Prev Page ∧K Cut Text ∧C Cur Po ∧X Exit ∧J Justify ∧W Where Is ∧Y Next Page ∧U UnCut Text ∧T To Spe	took only 12.35GB of RAM.

Pandas vs. Large Data Tips

- "Why and How to Use Pandas with Large (but not big) Data"
 - <u>https://towardsdatascience.com/why-and-how-to-use-pandas-with-large-data-9594dda2ea4c</u>
- Use generators to access files in chunks
- Filter out unimportant columns to make your DataFrame smaller
- Use optimal data types
 - Float64 is bigger than float32, but we might only need float32

Data Types and Optimization

- 1. List Comprehend through alice words list against enable words list.
 - [5]: %time notfound = [w for w in awords if w not in enable]

2. Same, but filter against enable words as a set.

[7]: %time notfound = [w for w in awords if w not in enable_set]

- 3. Compute set difference.
 - [9]: %time notfound = awords_set.difference(enable_set)

Data Types and Optimization

1. List Comprehend through alice words list against enable words list.

Ew

Best

option

- [5]: %time notfound = [w for w in awords if w not in enable] CPU times: user 39.7 s, sys: 16.2 ms, total: 39.7 s Wall time: 39.8 s
- 2. Same, but filter against enable words as a set.
 - [7]: %time notfound = [w for w in awords if w not in enable_set]
 CPU times: user 20.3 ms, sys: 4 ms, total: 24.3 ms
 Wall time: 24.8 ms
 faster
- 3. Compute set difference.
 - [9]: %time notfound = awords_set.difference(enable_set)
 CPU times: user 332 μs, sys: 2 μs, total: 334 μs
 Wall time: 339 μs

List objects are not optimized for membership operations, but sets are!

Algorithmic Efficiency: Summary

- Solutions can be implemented with varying degrees of efficiency
- Certain problems have inherent limits to efficiency
 - Big O represents the upper bound to a problem as it relates to processing time and the input size
 - Example: Sorting algorithms can only be as fast as O(n log n) in the worst case
 - There are some sorting algorithms that are faster, but they only work with special input types
- Take Away
 - Use the most efficient algorithm possible
 - "efficient" can be in terms of time or space complexity
 - Weigh your options and choose what's best for your situation
 - Understand the relationship between the growth function and the input size
 - O(n) is rather scalable, but O(n²) is too inefficient for large inputs
 - More efficient algorithms = better memory usage and faster runtimes
 - Efficiency is not always the be-all and end-all
 - With smaller data sets, maybe readability is more important