

Administrivia

Course overview

"Document distance" problem

"Peak finding" problem Intro

Handouts

Course-info

doc dist python code

Sign up for class at alg.csail.mit.edu
Read collaboration policy!

Course Overview

- Efficient procedures for solving problems on large inputs (e.g., US highway map, human genome)
- Scalability
- Classic data structures and elementary algorithms (CLRS text)
- Real implementations in Python
- Fun problem sets

Content

7 modules, each with motivating problem and problem set (except last)

- Linked data structures : Doc distance, peak finding
- Hashing : Selective file update (RSYNC)
- Sorting : gas simulation (?)
- Search : Rubik's cube (?)
- Shortest paths : CalTech → MIT (?)
- Dynamic Programming : Stock Market
- Numerics : $\sqrt{2}$ to infinite precision

Document Distance Problem

- Given two documents, how similar are they? identical is easy, plagiarized harder
- Need to define metric
- Word is sequence of alphanumeric characters
"6.006 is fun" 4 words
- Word frequencies : $D(w) = \# \text{ times } w \text{ occurs in document } D$
 count :

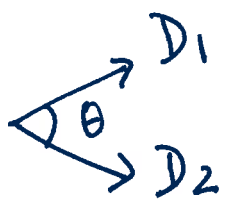
1	0	1	1	0	1
6	the	is	006	easy	fun

METRIC

$$D_1 \circ D_2 = \sum_w D_1(w) \cdot D_2(w)$$

inner product

$$\|D\| = N(D) = \sqrt{D \circ D}$$



$$\theta(D_1, D_2) = \arccos \left(\frac{D_1 \circ D_2}{\|D_1\| \cdot \|D_2\|} \right)$$

Identical $\rightarrow 0 \leq \theta \leq \pi/2$ \leftarrow no common words

PYTHON IMPLEMENTATION

docdist1.py

- Read file
- Make word list \rightarrow ["the", "year", ...]
- Count frequencies \rightarrow [{"the": 4012}, {"year": 55}, ...]
- Sort into order \rightarrow [{"a": 3120}, {"after": 17}, ...]
- Compute θ

Jules Verne	25K
Bobsey Twins	268K
Lewis & Clark	1M
Shakespeare	5.5M
Churchill	10M

Expt: Bobsey vs Lewis
 $\theta = 0.574$ ($> 2min$)

Dies on bigger files

What is going on?

Python vs C?
choice of algorithm?

Profiling

How much time spent in each routine

```
[ import profile
  profile.run("main()")
```

- ① # calls
- ② tot time : exclusive of subroutine calls
- ③ per call : ②/①
- ④ cum: including subroutine calls
- ⑤ per call : ④/①

Bobsey vs Lewis

Total: 141s
 get words from line list : 53s
 count-frequency: 50s
 get words from string : 12s
 insertion sort : 13s

BIGGEST CULPRIT

```
get_words_from_line_list(L):
  word_list = []
  for line in L:
    words_in_line = get_words_from_string(line)
    word_list = word_list + words_in_line
  return word_list
```

has to be this!
 (there isn't anything else here)

LIST CONCATENATION

(5)

$$L = L_1 + L_2$$

takes time proportional to $|L_1| + |L_2|$

Suppose we had n lines, each with one word

$$\text{time proportional to } 1 + 2 + 3 + \dots + n = \frac{n(n+1)}{2} = \theta(n^2)$$

Solution: `word-list.extend(words_in_line)`
i.e., `L1.extend(L2)`

This takes time proportional to $|L_2|$

for each word in `words_in_line` `word-list.append(word)`
get words from line list: 53s \rightarrow 0.12s

IMPROVEMENTS

docdist 1.py	original code	141s
docdist 2.py	add profiling	141s
docdist 3.py	<code>word-list.extend</code>	94s
docdist 4.py	dictionaries in count freq	42s
docdist 5.py	process words rather than chars in <code>get words from string</code>	17s
docdist 6.py	merge sort rather than insertion sort	6s
docdist 6B.py	Eliminate sorting by using dictionaries	0.5s

PEAK FINDER

One-dimensional version

1	2	3	4	5	6	7	8
a	b	c	d	e	f	g	h

a to h are numbers

Position 2 is a peak if and only if $b \geq a$ and $b \geq c$

Position 8 is a peak if $h \geq g$

Problem: Find a peak if it exists