

Administrivia

Course Overview

"Document distance" problem

Handouts

- #1 Course-info
- #2 Signup sheet
- #3 Doc distance

READ COLLABORATION POLICY!

PRE-REQS: Python, discrete mathematics will have 4 times that minimize conflicts

Course Overview

- Efficient procedures for solving problems on large inputs (entire works of Shakespeare, human genome, US highway map)
- Scalability
- Classic data structures and elementary algos (CLRS. text)
- Real implementations in Python/
Fun problem sets
- B version of class — give us feedback!

Content

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

- Linked data structures : Document distance (DD)
- Hashing : DD, genome comparison
- Sorting : gas simulation
- Search : Rubik's cube 2x2x2
- Shortest paths : Caltech → MIT
- Dynamic Programming : Stock market
- Numerics : $\sqrt{2}$

Document Distance Problem

- Given two documents, how similar are they?
 identical - easy?
 modified or related (DNA, plagiarism, authorship)
 Did Francis Bacon write Shakespeare's plays?
- Need to define metric
- Word is sequence of alphanumeric characters
 "6.006 is fun" 4 words.
- Word frequencies : $D(w)$ = #times w occurs in document D
 count : [1 0 1 1 0 1]
 W : 6 the is 006 easy fun

choose some canonical ordering of words

Profiling

How much time spent in each routine

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

① #calls

② tottime: exclusive of subroutine calls

③ per cell: ②/①

④ cum: including subroutine calls

⑤ percell: ④/①

Bobsey vs Lewis

- 195 secs total.
- 107: get words from line list
- 44: count-frequency
- 13: get words from string
- 12: insertionsort

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

$$L = L_1 + L_2$$

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

if we had n lines, each with one word

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

Solution: word-list.extend(words-in-line)

$L_1.extend(L_2)$ time proportional to $|L_2|$ [word-list.append(word)] for each word in words-in-line

Python has powerful primitives built in. To write efficient algorithms, we need to understand their costs. Figure out cost of set operations in PS 1 ...

docdist3.py

195 secs \rightarrow 85 secs.

```

def count-frequency (word-list)
  L = []
  for new-word in word-list :
    for entry in L :
      if entry[0] == new-word:
        entry[1] = entry[1] + 1
        break
    else:
      L.append([new-word, 1])
  return L

```

Analysis:

n words
 d distinct words
 $\theta(nd)$
if all words distinct $\theta(n^2)$

Dictionaries

Hash tables: mapping from domain
 (finite collection of immutable things) to
 range (anything). $D['ab'] = 2$ $D['the'] = 3$

docdist4.py uses dictionaries that give
 constant time lookup. 85 secs \rightarrow 42 secs

see count-frequency. $D = \{ \}$
 for new-word in word-list:
 if D.has-key(new-word):
 $D[new-word] = D[new-word] + 1$
 else:
 $D[new-word] = 1$
 return D.items()

time $\Theta(n)$

What's left? get-words-from-string 13 secs
 (V5 fixes with translate)
 insertion-sort 11 secs
 (V6 fixes with merge-sort)

next time