

(turn cellphone off!)

6.006

Rivest

Outline:

L1.1

9/4/08

- Administrivia
- Course overview
- "Document distance" problem

Handouts: 1. Course info

2. Sign-up sheet (passed around)
3. docdist1.py

Administrivia:

- Welcome to 6.006!

- Introduce staff.

- Sign-up - on line

and - on sheet being passed around

{ web = <http://courses.csail.mit.edu/6.006>

- Ask students to raise hands for:

credit/listener/unsure

frash/soph/...

MIT/not

course 6 / 18 other

python

6.01

6.042

- Pre-regs: see TA's if haven't got them, but have \equiv
- Recitations: on-line; get assignment tonight - bring laptops!
- lectures/recitations/problemsets / 2 quizzes / final
- textbook : CLRS, rec: Miller/Ranum
- relation to 6.046

Course Overview

- Efficient procedures for solving problems on large inputs
- scalability (now can have, & we'll use complete works of Shakespeare, human DNA, or U.S. highway map on your laptop)
- classic data structures & elementary algorithms (CLRS)
- real implementations (Python)
- having some fun (problem sets!)
- developing this course:
 - this is only ~~initial~~ version - will have rough edges
 - we want your feedback - think of yourselves as "co-designers"

Content:

- 7 modules, each with motivating problem & problem set (except last)
- Intro & linked data structures: Document Distance Set Ops
- Hashing " "
- Dynamic Programming Image Resizing → ~~Sequence Alignment~~
 - Sorting Chimp DNA
 - Search Gas Simulation
 - Shortest Paths 2x2x2 Rubik
 - Numerics Cal Tech → MIT
- Least Squares

6,006
 Rivest
 L1.3
~~1~~
 9/4/08

Document Distance Problem (Document Similarity)

- Given two "documents" how similar are they? common problem
- identical - easy?
- modified or related (DNA, plagiarism, authorship)
- Did Francis Bacon write Shakespeare's plays?
- Need to define metric
- define word = sequence of alphanumerics "6,006 is fun" 4 words

- word frequencies: $D(w) = \#$ times w occurs in document D

can think of frequencies as a vector, each word is a coordinate

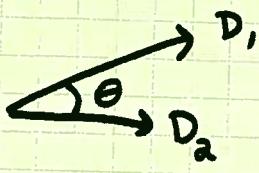
Count $[1, 0, 1, 1, 0, 1]$

w 6, the, is, 006, easy, fun

for some canonical ordering of words...

$$D_1 \cdot D_2 = \sum_w D_1(w) \cdot D_2(w) \quad \text{inner product}$$

$$\|D\| = N(D) = \sqrt{D \cdot D} \quad \text{length, norm}$$



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

$$0 \leq \Theta \leq \pi/2$$

↑
identical

↑
no common
words

Problem: given D_1, D_2 compute Θ (angle between their word frequency vectors)

	<u>(bytes)</u>	6,006 L1.4 Rivest 9/4/08
• Data Set 3:		
Jules Verne "2889"	25K	
Bobsey Twins	268K	
Lewis & Clark	1M	
Arabian Nights	3M	
Churchill	10M	
Shakespeare	5.5M	
Bacon	320K	

Project Gutenberg

(You can also try
reading these!)

(esp. Verne) → read
excerpt

Procedure (docdist1):

- read file
- make word list $["\text{the}", "\text{year}", "2889", "by", "Jules", \dots]$
- $\times 2$ • Count frequencies $[[\text{the}, 4012], [\text{year}, 55], \dots]$
- Sort into order $[[\text{a}, 3120], [\text{after}, 17], \dots]$
- Compute $\theta = \arccos\left(\frac{\mathbf{D}_1 \cdot \mathbf{D}_2}{\|\mathbf{D}_1\| \cdot \|\mathbf{D}_2\|}\right)$

Look at Python code:

- students should be able to read this code (ask TA if not, in recit.)
- go through routines (admittedly not most efficient, but correct.)

Experiment: Bobsey vs Lewis $\theta = 0.574$ (3 minutes)

Doesn't seem to scale well, just "dies" on bigger files... (too long)

What is going on?

DSSS [Python vs C? (choice of programming language) $\times 10$ or so only
choice of algorithm $\Theta(n^2) \Rightarrow \Theta(n)$ 10^3 or more speedup
 \sqrt{n} much more important

6,006
L1,5
River
9/4/08

Profiling:

- How much time is spent in each routine:
~~total execution of subroutine~~
~~parent~~

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

- ① # calls
- ② tottime - exclusive of subroutine calls
- ③ percall - ②/①
- ④ cum - including subroutine calls
- ⑤ percum - ④/①

Results: Bob vs Lewis

- 194 seconds total - 3 minutes!
- 107 in get_words from line list
- 44 in count_frequency
- 13 in get_words from string
- 12 in insertion-sort

} exclusive of
subroutine calls
they make

} eventually
we'll get
all this
down to
6 seconds

[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!)

6.006
Rivest
~~L1.6~~
9/4/08

List Concatenation:

$$L = L_1 + L_2$$

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

if we had n lines, each with ~~one~~ 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)` \leftarrow time proportional to $\text{len}(\text{words-in-line})$

`[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.

(Figuring out cost of set opns will be their homework...)

<u>Document Distance</u>	t2.bobsey.txt t3.lewis.txt	263 KB 1 MB	6,006 Rivest L1.7 9/4/08 ? secs ✓
Version 1: initial			
2: add profiling			194 secs ✓
3: wordlist.extend(words-in-line)			84 secs ✓
4: use dictionaries in count-frequency			41 secs
5: translate			13 secs
6: merge-sort			6 secs

with docdist3:

def count-frequency(word-list):

L = []

for new-word in word-list:

 for entry in L:

 if entry[0] == w:

 entry[1] = entry[1] + 1

 break

 else:

 L.append([new-word, 1])

return L

try:

 from docdist3 import *

 L = read-file("t2.bobsey.txt")

L[0], L[-7]

 W = get-words-from-line-list(L)

W[:7], W[-7:]

 F = count-frequency(W)

(slow on lewis - don't try in class)

Analysis: time = $O(n \cdot d)$

if all words distinct, $d = n$

n words

d distinct words

time = $\Theta(n^2)$

Dictionaries - Hash Tables

mapping from domain (finite collection of immutable things)
to range (anything)

$D = \{\}$ empty mapping

$D['ab'] = 2$

$D['the'] = 3$ fast!

D

$D['ab']$ fast!

$D['xyz']$ error

`D.has_key('xyz')`

`D.items()`

`D.keys()`

Show `docdist4 count_frequency`

cuts time in $1/2$ (more for larger files) ($84\text{ secs} \Rightarrow 41\text{ secs}$)

Analysis: time $\Theta(n)$ doesn't depend on what's already in table

Remaining time goes to:

`get_words_from_string` (vs `translate`) 13 secs

`insertion-sort` (vs `merge-sort`) 11 secs

Important to understand costs (running times) of Python primitives!

See Python Cost Model (web site) for some experimentation...

(You'll mimic this for Python `set` operations...) HW #1