

(turn cellphone off!)

6.006

original: Rivest

Outline:

L1.1

9/4/08

□ Administrivia

□ Course overview

update: Edelman

□ "Document distance" problem

2/3/09

Handouts: 1. Course info

2. sign-up sheet (passed around)

3. docdist1.py

Administrivia:

• Welcome to 6.006!

• Introduce staff

• sign-up - on line (see announcements)

and - on sheet being passed around

• web = <http://courses.csail.mit.edu/6.006> (Spring 2009)

• Ask students to raise hands for:

credit/listener/unsure

fresh/soph/...

MIT/not

course 6 / 18 other

python

6.01

6.042

• Pre-regs: see TA's if haven't got them, but have ≡

• Recitations: on-line; get <sup>PS</sup> assignment tonight - ~~assignments~~

• lectures/recitations/pubblmscts / 2 quizes / final (3/11, 4/15)

• text book: CLRS, rec: Miller/Reagan

34-401 50-340

• relation to 6.046

## Course Overview

- Efficient procedures for solving problems on large inputs
- scalability (now can have, & we'll use complete words 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: (AE: Now iteration #4!)
  - this is only ~~first~~ version - will have rough edges
  - we want your feedback - think of yourselves as "co-designers"

6.006  
Rivest  
L1.2  
~~9/4/08~~  
AE: 2/3/09

## Sample Content:

- 7 modules, each with motivating problem & problem set (except last)
  - Intro & linked data structures: Document Distance      Set Opns
  - Hashing      "      ~~Chimp DNA~~
  - Dynamic Programming: Image Resizing → ~~Chimp DNA~~
  - Sorting      Heapsort      Gas Simulation
  - Search      BF/DFS      2x2x2 Rubik
  - Shortest P-ths      CalTech → MIT
  - Numics      ~~Least Squares~~  
Least Squares

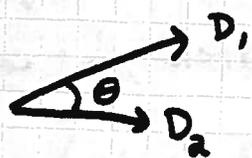
## Document Distance Problem (Document Similarity)

6.006  
Rivest  
L1.3  
~~AE~~  
9/4/08  
AE/2/3/09

- 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 <sup>possible</sup> word is a coordinate
- Count  $[1, 0, 1, 1, 0, 1]$   
 $w$  6, the, is, 006, easy, fun  
for some canonical ordering of words...

$$\mathbf{D}_1 \cdot \mathbf{D}_2 = \sum_w D_1(w) \cdot D_2(w) \quad \text{inner product}$$

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



$$\theta = \arccos\left(\frac{\mathbf{D}_1 \cdot \mathbf{D}_2}{\|\mathbf{D}_1\| \cdot \|\mathbf{D}_2\|}\right) = \theta(\mathbf{D}_1, \mathbf{D}_2) \quad \text{correlation}$$

$$0 \leq \theta \leq \pi/2$$

↑ identical                      ↑ no common words

AE: or no arccos  
giving  $p = \cos \theta$   
 $0 \leq p \leq 1$

Problem: given  $\mathbf{D}_1, \mathbf{D}_2$  compute  $\theta$  (angle between their word frequency vectors)

Open question if this distance is a "good" measure or not. You might think why/why not

• Data Sets:		(bytes) ← AE's Numbers Vary	6,006 L1.4 <del>Rivest</del> 9/4/08 Elema 2/3/01
	Jules Verne "2889"	25K	Project Gutenberg (You can also try reading these!) (esp. Verne) → read excerpt AE: Good Idea!
	Bobby Twins	268K	
	Lewis & Clark	1M	
	Arabian Nights	3M	
	Churchill	10M	
	Shakespeare	5.5M	
	Bacon	320K	

Procedure (docdist1):

- read file
- make word list ["the", "year", "2889", "by", "Jules", ...]
- x2 ↑ two files input • count frequencies [ ["the", 4012], ["year", 55], ... ]
- sort into order [ ["a", 3120], ["after", 17], ... ]
- compute angle =  $\arccos\left(\frac{D_1 \cdot D_2}{\|D_1\| \cdot \|D_2\|}\right)$

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?

Discus [ Python vs C? (choice of programming language)  $\times 10$  or so only (AE's?)  
choice of algorithm  $\Theta(n^2) \Rightarrow \Theta(n)$   $10^3$  or more speedup  
 much more important AE: these days complicated software/hardware issues also relevant but not focus of this class

## Profiling:

- How much time is spent in each routine:  
~~total - exclusive of subroutines~~

import profile  
profile.run("main()") Edelman  
2/3/09

6.006  
L1.5  
Rivest  
~~9/4/08~~  
9/4/08

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

## Results: Bobers 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 inbuilt-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!)

## 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)` ← 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 ops will be their homework...)

6.006  
Rivest  
Li, G  
~~1/17~~  
9/4/08  
Edelman  
2/3/09

<u>Document Distance</u>	t2. bobsey.txt t3. lewis.txt	263 KB 1 MB	Edelman 2/3/09	6,006 Rivest <del>1.7</del> L1.7 2/4/08
Version 1 :	initial			? secs ✓
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 = []
```

```
[["the", 2509], ["on", 369]]
```

```
    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 index)
```

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

if all words distinct,  $d = n$

time =  $\Theta(n^2)$

n words  
d distinct words

## Dictionaries - Hash Tables

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

6.006  
Rivest  
~~9/4/08~~ L1.8  
9/4/08  
Edelman  
2/3/09

$D = \{ \}$  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 doctest4 count\_frequency

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

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

Remaining time goes to:

get\_words\_from\_string (VS times with translate) 13 secs

insertion\_sort (VS times with 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