_ILDER_
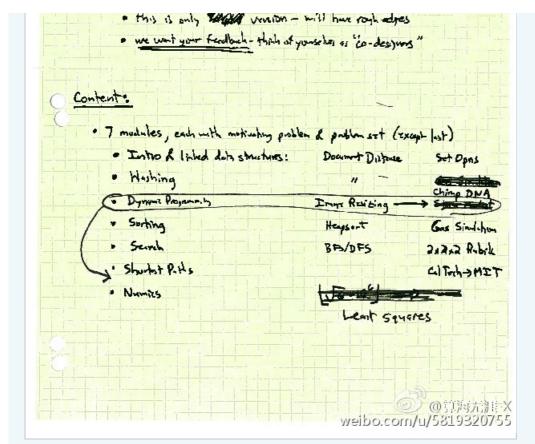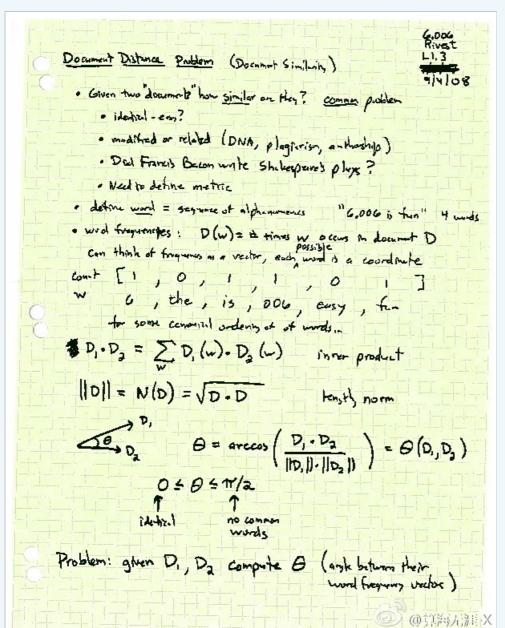
@算海无涯-X

#CS名校算法资源# MIT 6.006算法课考古贴: http://t.cn/RqdvHYu 收录了从2007年至今的课程资源，有两个发现：一是他们一直使用手写讲义并且传承悠久；二是从2007年开始就已经使用Python作为描述语言。附图是RSA密码算法字母R的那位大神Rivest的2008年珍贵手写讲义摘选： http://t.cn/RqdvHYn

- this is only ~~rough~~ 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 Opns
  - Hashing                                    "        "        ~~████~~
                                                                China DNA
  - Dynamic Programming          Image Resizing ⟶        ~~████~~
  - Sorting                              Heapsort          Gas Simulation
  - Search                               BFS/DFS          2x2x2 Rubik
  - Shortest Paths                                         CalTech→MIT
  - Numerics                             ~~████~~

                                           Least Squares

---

## 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 possible word is a coordinate

  count $\begin{bmatrix} 1 & , & 0 & , & 1 & , & 1 & , & 0 & 1 \end{bmatrix}$
  w      6   ,  the ,  is ,  006 ,  easy ,  fun
     for some canonical ordering of of words...

  $D_1 \cdot D_2 = \sum_w D_1(w) \cdot D_2(w)$     inner product

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



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

  $0 \le \theta \le \pi/2$
      ↑            ↑
   identical   no common
                words

Problem: given $D_1, D_2$ compute $\theta$ (angle between their
                                      word frequency vector)