

Lecture 2: Models of Computation, Document Distance

□ Algorithm

al-khwārizmī: Father of algebra

◦ What's an algorithm

- Computational procedure for solving a problem

Input \rightarrow alg. \rightarrow output

program \leftrightarrow algorithm

programming language

 \leftrightarrow

pseudo code structured English

Computer \leftrightarrow model of
computation

◦ Model of Computation

- what operation on algorithm is allowed
- Cost (time, ...) of each operation

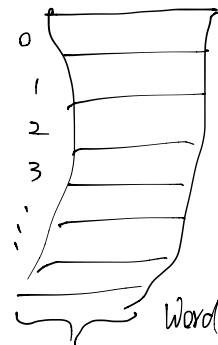
• Models

① Random Access Machine (RAM)

- Random access memory

modeled by big array

- In const. time an algorithm can
read in/ load a const. number of



words from memory, do a const. num of operations on them.
and store (write out) them.

- Const. registers.

- A word is w bits, $w \geq \lg(\text{size of memory})$

② Pointer Machine

- Dynamically allocated objs.

- obj has const. num of fields

- field = word (e.g. int) or pointer to obj. or null (None)

□ Python Model

① "list" = array

$L[i] = L[i] + 5$ takes const. time

② obj. with const. num of attributes ("pointer machine model")

$X = X.net$ takes const. time

$L.append(X)$ python does stable doubling [lecture 9]

can be done in const. time

$L = L1 + L2$ Take linear time

$\equiv L = []$

$\left. \begin{array}{l} \text{for } x \text{ in } L1: \underbrace{L.append(x)}_{\text{const. time}} \\ \text{for } x \text{ in } L2: \underbrace{L.append(x)}_{\text{const. time}} \end{array} \right\} \begin{array}{l} \Theta(|L1|) \\ \Theta(|L2|) \end{array}$

$\Theta(1 + |L1| + |L2|)$

$x \text{ in } L$: Take linear time

$\text{len}(L)$: Const time, Lists in python are implemented with the length of the list build in. \rightarrow Instantaneous

$L.\text{sort}()$ take $\Theta(n \log n)$ time

- time to compare items
 - python uses algorithm "Comparison Sort"
- (lecture 3)

dict : $D[\text{key}] = \text{val.}$

- Takes const. time with high probability (randomized algorithm)
- Hash table (lecture 8 - 10)

long : x, y has $|x|, |y|$ words long

- $x + y$ takes $\Theta(|x| + |y|)$
- $x \cdot y$ takes $\Theta((|x| + |y|)^{\log_2 3})$, while straight forward method takes $\Theta(|x| \cdot |y|)$ time

heapq : lecture 4

□ Document Distance Problem

Given two documents $\mathcal{D}_1, \mathcal{D}_2$, Compute the distance between them

• Distance: Are they identical

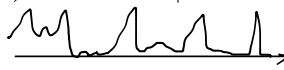
- document: sequence of words

- word: string of alphanumeric chars

- idea: look at shared words, use that to define distance

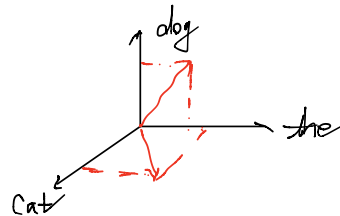
- Think of document as a vector

$\mathcal{D}[W] = \#$ of occurrences of W in \mathcal{D}

 words

$\mathcal{D}_1 = \text{"the cat"}$

$\mathcal{D}_2 = \text{"the dog"}$



Document distance \longleftrightarrow Vector distance

$$d(\mathcal{D}_1, \mathcal{D}_2) = \arccos\left(\frac{\mathcal{D}_1 \cdot \mathcal{D}_2}{|\mathcal{D}_1| |\mathcal{D}_2|}\right)$$

Algorithm:

- ① Split document into words
- ② Compute word frequencies
- ③ dot product

Implement in python:

②: iter thru doc.

$$\text{count[word]} += \frac{1}{\text{H}(|\text{word}|)}$$

"H(|word|)"

① Scan ~~thru~~ and look for alpha numerics