

Objectives

- Wrap up Data Compression/Huffman Codes
- Divide and conquer
 - Recurrence relations

Mar 1, 2010

CSCI211 - Sprenkle

1

Review: Prefix Codes

- **Problem:** Encoding of one character is a *prefix* of encoding of another
- **Solution: Prefix Codes:** map letters to bit strings such that *no encoding is a prefix of any other*
 - Won't need artificial devices like spaces to separate characters
- **Example encodings:**
 - Verify that no encoding is a prefix of another
 - What is 0010000011101?

a: 11	d: 10
b: 01	e: 000
c: 001	

Mar 1, 2010

CSCI211 - Sprenkle

2

Review: Optimal Prefix Codes

- **Goal:** minimize **Average number of Bits per Letter (ABL):**

$$\sum_{x \in S} \text{frequency of } x * \text{length of encoding of } x$$

↑ For all characters in our alphabet
- f_x : frequency that letter x occurs
- $\gamma(x)$: encoding of x
 - $|\gamma(x)|$: length of encoding of x
- Minimize **ABL** = $\sum_{x \in S} f_x |\gamma(x)|$

Mar 1, 2010

CSCI211 - Sprenkle

3

Review: Problem Statement

- Given an alphabet and a set of frequencies for the letters, produce optimal (most efficient) prefix code
 - Minimizes average number of bits per letter

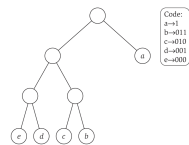
Mar 1, 2010

CSCI211 - Sprenkle

4

Review: Binary Trees to Represent Prefix Codes

- Exposes structure better than list of mappings
 - Each leaf node is a letter
 - Follow path to the letter
 - Going left: 0
 - Going right: 1



Mar 1, 2010

CSCI211 - Sprenkle

5

Review: Huffman Algorithm

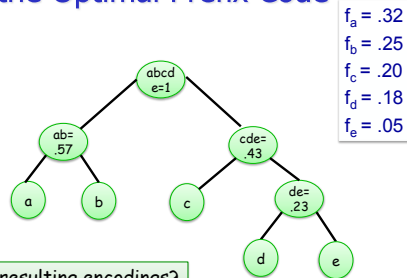
1. Create a leaf node for each symbol, labeled by its frequency, and add to a queue
2. While there is more than one node in the queue
 - a) Remove the two nodes of lowest frequency
 - b) Create a new internal node with these two nodes as children and with frequency equal to the sum of the two nodes' probabilities
 - c) Add the new node to the queue
3. The remaining node is the tree's root node

Mar 1, 2010

CSCI211 - Sprenkle

6

Example: Creating the Optimal Prefix Code



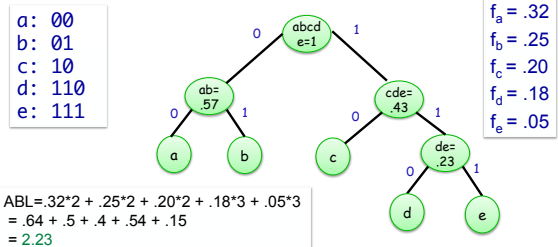
What are the resulting encodings?
What is the ABL?

Mar 1, 2010

CSCI211 - Sprenkle

7

Creating the Optimal Prefix Code



I chose to build the tree this way.
What if I had switched the order of the children?

Mar 1, 2010

CSCI211 - Sprenkle

8

Implementation

- What data structures do we need?

Mar 1, 2010

CSCI211 - Sprenkle

9

Implementation

- What data structures do we need?
 - Binary tree for the prefix codes
 - Priority queue for choosing the node with lowest frequency
- Where are the costs?

Mar 1, 2010

CSCI211 - Sprenkle

10

Running Time

- Costs
 - Inserting and extracting node into PQ: $O(\log n)$
 - Number of insertions and extractions: $O(n)$
 - $O(n \log n)$

Mar 1, 2010

CSCI211 - Sprenkle

11

Analysis of Algorithm's Optimality

- 2 page proof in book

Mar 1, 2010

CSCI211 - Sprenkle

12

Real-life Compression

- Text can be compressed well because of known frequencies
- Algorithms can be optimized to languages
 - More than just “z doesn’t happen very often”
 - “z doesn’t happen after q”

Mar 1, 2010

CSCI211 - Sprenkle

13

DIVIDE AND CONQUER ALGORITHMS

Mar 1, 2010

CSCI211 - Sprenkle

14

Divide-and-Conquer

Divide et impera.
Veni, vidi, vici.
- Julius Caesar

- Divide-and-conquer process
 - **Break up** problem into **several parts**
 - Solve each part **recursively**
 - **Combine** solutions to sub-problems into overall solution
- Most common usage:
 - Break up problem of size n into two equal parts of size $\frac{1}{2}n$
 - Solve two parts recursively
 - Combine two solutions into overall solution

Mar 1, 2010

CSCI211 - Sprenkle

15

Discussion

- What is a well-known divide and conquer algorithm?

MERGE SORT

Mar 1, 2010

CSCI211 - Sprenkle

16

Merge Sort

- How does Merge Sort work?
- When do we stop?

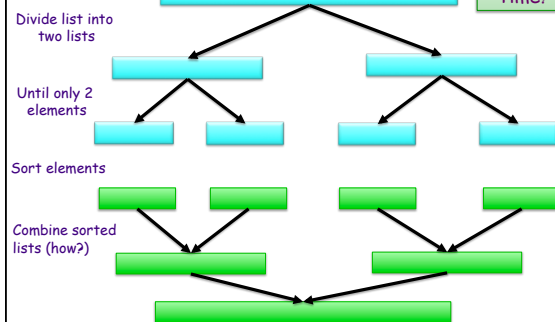
Mar 1, 2010

CSCI211 - Sprenkle

17

Merge Sort

Costs?
Running
Time?



Mar 1, 2010

CSCI211 - Sprenkle

18

RECURRENCE RELATIONS

Mar 1, 2010

CSCI211 - Sprenkle

19

Analyzing Merge Sort

General Template

- Break up problem of size n into two equal parts of size $\frac{1}{2}n$
- Solve two parts recursively
- Combine two solutions into overall solution

- **Def.** $T(n)$ = number of comparisons to mergesort an input of size n
- Want to say a bit more about what $T(n)$ is
 - Break it down more...
 - What can we say about the running time w.r.t. to the different parts of the above template?

Mar 1, 2010

CSCI211 - Sprenkle

20

Analyzing Merge Sort

General Template

- Break up problem of size n into two equal parts of size $\frac{1}{2}n$ $O(1)$
- Solve two parts recursively $T(n/2) + T(n/2)$
- Combine two solutions into overall solution $O(n)$

- **Def.** $T(n)$ = number of comparisons to mergesort an input of size n
- Want to say a bit more about what $T(n)$ is
 - Break it down more...
 - What can we say about the running time w.r.t. to the different parts of the above template?

Mar 1, 2010

What is the base case's running time?

Merge Sort's Recurrence Relation

- Put an *upperbound* on $T(n)$:

For some constant c ,
 $T(n) \leq 2T(n/2) + cn$ when $n > 2$,
 $T(2) \leq c$.

Why is this constant?

Solve $T(n)$ to come up with explicit bound

Mar 1, 2010

CSCI211 - Sprenkle

22

Approaches to Solving Recurrences

1. Unroll recursion

- Look for patterns in runtime at each level
- Sum up running times over all levels

2. Substitute guess solution into recurrence

- Check that it works
- Induction on n

Mar 1, 2010

CSCI211 - Sprenkle

23

Unrolling Recurrence: $T(n)$

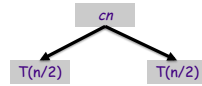
Mar 1, 2010

CSCI211 - Sprenkle

24

Unrolling Recurrence

- First level: $2 T(n/2) + cn$



How does the next level break down?

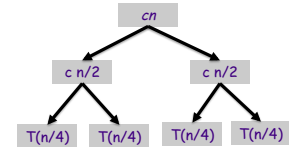
Mar 1, 2010

CSCI211 - Sprenkle

25

Unrolling Recurrence

- Next level:



Each one is $2 T(n/4) + c(n/2)$

Next level?

Mar 1, 2010

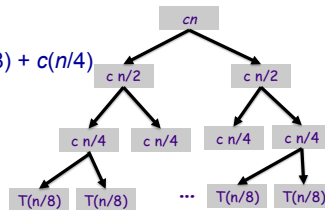
CSCI211 - Sprenkle

26

Unrolling Recurrence

- Next level:

Each one is $2 T(n/8) + c(n/4)$



And so on...

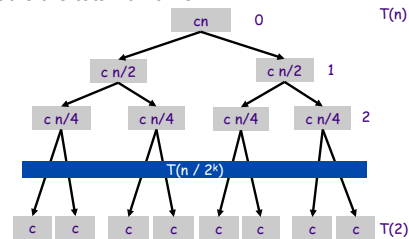
Mar 1, 2010

CSCI211 - Sprenkle

27

Unrolling Recurrence

- How much does each level cost, in terms of the level?
- How many levels are there (assuming n is a power of 2)?
- What is the total run time?



Mar 1, 2010

CSCI211 - Sprenkle

28

Unrolling Recurrence

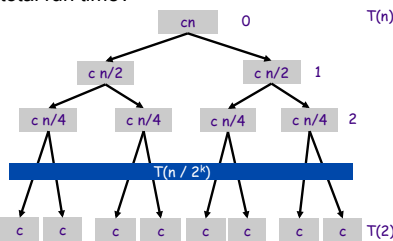
- How many levels are there (assuming n is a power of 2)?
- How much does each level cost, in terms of the level?
- What is the total run time?

Number of levels:
 $\log_2 n$

2^k problems
Size: $n/2^k$

Each level takes
 $2^k * c * (n/2^k) = cn$

$\Rightarrow O(n \log n)$



Mar 1, 2010

CSCI211 - Sprenkle

29

Alternative: Proof by Induction

- Claim.** If $T(n)$ satisfies this recurrence, then $T(n) = n \log_2 n$.

$$T(n) = \begin{cases} 0 & \text{if } n = 1 \\ 2T(n/2) + n & \text{otherwise} \end{cases}$$

(sorting both halves merging)

- Pf.** (by induction on n)

➤ Base case: $n = 1$

➤ Inductive hypothesis: $T(n) = n \log_2 n$

➤ Goal: show that $T(2n) = 2n \log_2 (2n)$

Why doubling n ?

Mar 1, 2010

CSCI211 - Sprenkle

30

Proof by Induction

- **Claim.** If $T(n)$ satisfies this recurrence, then $T(n) = n \log_2 n$.

$$T(n) = \begin{cases} 0 & \text{if } n=1 \\ \underbrace{2T(n/2)}_{\text{sorting both halves}} + \underbrace{n}_{\text{merging}} & \text{otherwise} \end{cases}$$

- **Pf.** (by induction on n)
 - Inductive hypothesis: $T(n) = n \log_2 n$

$$\begin{aligned} T(2n) &= 2T(n) + 2n \\ &= 2n \log_2 n + 2n \\ &= 2n(\log_2(2n) - 1) + 2n \\ &= 2n \log_2(2n) \end{aligned}$$

Mar 1, 2010

CSCI211 - Sprenkle

31

Another Example

- Instead of recursively solving 2 problems, solve q problems
 - Size of problems is still $n/2$
- Combining solutions is still $O(n)$

Mar 1, 2010

CSCI211 - Sprenkle

32

Another Example

- Instead of recursively solving 2 problems, solve q problems
 - Size of problems is still $n/2$
- Combining solutions is still $O(n)$
- **Recurrence relation:**
 - For some constant c ,

$$T(n) \leq q T(n/2) + cn \text{ when } n > 2$$

$$T(2) \leq c$$

Intuition about running time?

Mar 1, 2010

CSCI211 - Sprenkle

33

Assignments

- Continue reading Chapter 4
- Start reading Chapter 5
- Wiki
 - Read 4.6-4.8, 5.1
- PS5 due Friday

Mar 1, 2010

CSCI211 - Sprenkle

34