

Objectives

- Wrap up Huffman Codes
- Divide and Conquer Algorithms

Feb 27, 2013

CSCI211 - Sprenkle

1

Review

- What was the problem we were trying to solve on Monday?
 - What was our optimization goal?
- What was our solution to the problem?

Feb 27, 2013

CSCI211 - Sprenkle

2

Review: Problem and Goal

- Do we need an average of 5 bits/character always?
- What if we could use *shorter encodings* for *frequently* used characters, like a, e, s, t?

Goal: Optimal encoding that takes advantage of *nonuniformity* of letter frequencies

- A fundamental problem for **data compression**
 - Represent data *as compactly as possible*

Feb 27, 2013

CSCI211 - Sprenkle

3

Review: Combining Our Conclusions

- The binary tree corresponding to the optimal prefix code is *full*, i.e., each internal node has two children
- We want to label the leaf nodes of the binary tree corresponding to the optimal prefix code such that nodes with *greatest depth* have *least frequency*

What does this mean the bottom of our tree looks like?

Feb 27, 2013

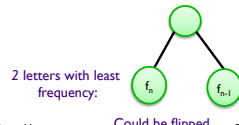
CSCI211 - Sprenkle

4

Review: Combining Our Conclusions

- The binary tree corresponding to the optimal prefix code is *full*, i.e., each internal node has two children
- We want to label the leaf nodes of the binary tree corresponding to the optimal prefix code such that nodes with *greatest depth* have *least frequency*

What does this mean the bottom of our tree looks like?



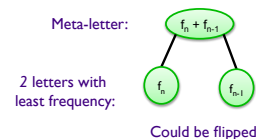
Feb 27, 2013

CSCI211 - Sprenkle

5

Review: How Can We Use This?

- Two letters with least frequency are definitely going to be siblings
 - Tie them together
 - Their parent is a “meta-letter”
 - Frequency is sum of $f_n + f_{n-1}$



Feb 27, 2013

CSCI211 - Sprenkle

6

Review: Huffman Coding Algorithm

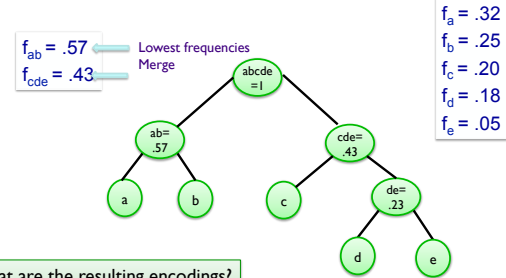
1. Create a leaf node for each symbol, labeled by its frequency, and add to a queue
2. While there is more than one node in the queue
 - a) Remove the two nodes of lowest frequency
 - b) Create a new internal node with these two nodes as children and with frequency equal to the sum of the two nodes'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

Feb 27, 2013

CSCI211 - Sprenkle

7

Creating the Optimal Prefix Code



What are the resulting encodings?
What is the ABL?

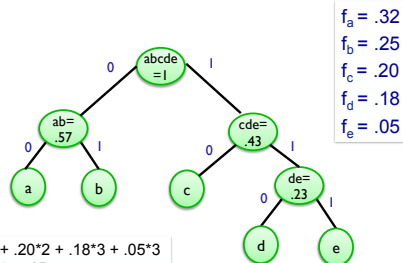
Feb 27, 2013

CSCI211 - Sprenkle

8

Creating the Optimal Prefix Code

$a: 00$
 $b: 01$
 $c: 10$
 $d: 110$
 $e: 111$



$$\begin{aligned}
 ABL &= .32 \cdot 2 + .25 \cdot 2 + .20 \cdot 2 + .18 \cdot 3 + .05 \cdot 3 \\
 &= .64 + .5 + .4 + .54 + .15 \\
 &= 2.23
 \end{aligned}$$

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

Feb 27, 2013

Implementing Huffman Coding 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

What data structures do we need?

- 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

Feb 27, 2013

CSCI211 - Sprenkle

10

Implementation

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

Feb 27, 2013

CSCI211 - Sprenkle

11

Analyzing Runtime of Huffman Coding Algorithm

Where are the costs?

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

Feb 27, 2013

CSCI211 - Sprenkle

12

Running Time

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

Feb 27, 2013

CSCI211 - Sprenkle

13

Analysis of Algorithm's Optimality

- 2 page proof in book

Feb 27, 2013

CSCI211 - Sprenkle

14

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"

Feb 27, 2013

CSCI211 - Sprenkle

15

DIVIDE AND CONQUER ALGORITHMS

Feb 27, 2013

CSCI211 - Sprenkle

16

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

Feb 27, 2013

CSCI211 - Sprenkle

17

Discussion

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

Merge Sort

Feb 27, 2013

CSCI211 - Sprenkle

18

Merge Sort

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

Feb 27, 2013

CSCI211 - Sprenkle

19

Merge Sort

Divide list
into two lists

Until only 2
elements

Sort elements

Combine sorted
lists (how?)

Feb 27, 2013

CSCI211 - Sprenkle

20

RECURRENCE RELATIONS

Feb 27, 2013

CSCI211 - Sprenkle

21

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?

Feb 27, 2013

22

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 is the base case? Its running time?

Feb 27, 2013

CSCI211 - Sprenkle

23

Merge Sort's Recurrence Relation

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

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

$O(n)$

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

Feb 27, 2013

CSCI211 - Sprenkle

24

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

Feb 27, 2013

CSCI211 - Sprenkle

25

Unrolling Recurrence: $T(n)$

$$T(n) = 2 T(n/2) + cn$$

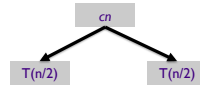
Feb 27, 2013

CSCI211 - Sprenkle

26

Unrolling Recurrence: $2 T(n/2) + cn$

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



How does the next level break down?

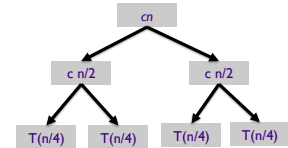
Feb 27, 2013

CSCI211 - Sprenkle

27

Unrolling Recurrence: $2 T(n/2) + cn$

- Next level:



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

Next level?

Feb 27, 2013

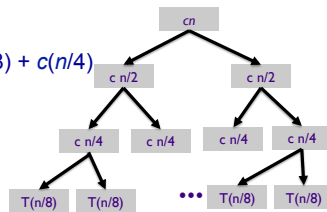
CSCI211 - Sprenkle

28

Unrolling Recurrence

- Next level:

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



And so on...

What does the final level look like?

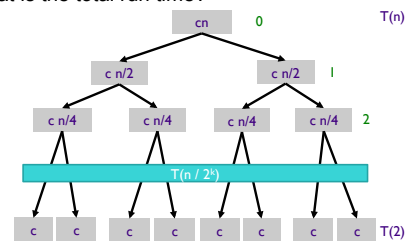
Feb 27, 2013

CSCI211 - Sprenkle

29

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?



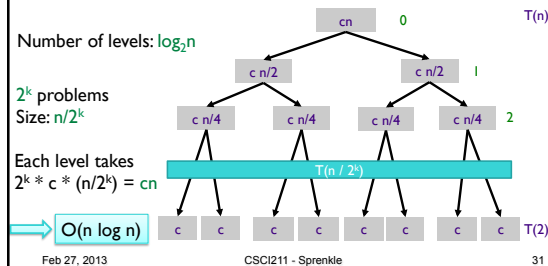
Feb 27, 2013

CSCI211 - Sprenkle

30

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?



Alternative: Proof by Induction

- Claim.** If $T(n)$ satisfies this recurrence, then $T(n) = n \log_2 n$.
 - Recall: $T(n) = 2 T(n/2) + cn$
- Pf.** (by induction on n)
 - Base case: $n = 2$
 - Inductive hypothesis: $T(n) \leq cn \log_2 n$
 - Goal: show that $T(2n) = 2cn \log_2 (2n)$

Why doubling n ?

Feb 27, 2013

CSCI211 - Sprenkle

32

Proof by Induction

- Claim.** If $T(n)$ satisfies this recurrence, then $T(n) = n \log_2 n$.
 - Recall: $T(n) = 2 T(n/2) + cn$
- Pf.** (by induction on n)
 - Inductive hypothesis: $T(n) \leq cn \log_2 n$
 - Goal: show that $T(2n) = 2cn \log_2 (2n)$

$$\begin{aligned}
 T(2n) &= 2T(n) + c2n \\
 &= 2cn \log_2 n + 2cn \quad \text{Replace w/ induction hypothesis} \\
 &= 2cn (\log_2(2n) - 1) + 2cn \\
 &= 2cn \log_2(2n) - 2cn + 2cn \\
 &= 2cn \log_2(2n) \quad \checkmark
 \end{aligned}$$

Feb 27, 2013

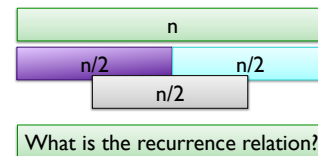
CSCI211 - Sprenkle

33

Another Recurrence Relation

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

Example: $q=3$:



Feb 27, 2013

CSCI211 - Sprenkle

34

Another Recurrence Relation

- 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?

Feb 27, 2013

CSCI211 - Sprenkle

35

Looking Ahead

- Problem Set 5 due Friday

Feb 27, 2013

CSCI211 - Sprenkle

36